

AD-A213 777

LABORATORY FOR  
COMPUTER SCIENCE



MASSACHUSETTS  
INSTITUTE OF  
TECHNOLOGY

FILE COPY

4

MIT/LCS/TM-406

# ON THE CORRECTNESS OF ORPHAN MANAGEMENT ALGORITHMS

Maurice Herlihy  
Nancy Lynch  
Michael Merritt  
William Weihl

DTIC  
ELECTE  
OCT 30 1989  
S E D

This document has been approved  
for public release and sale in  
distribution is unlimited.

August 1989

545 TECHNOLOGY SQUARE, CAMBRIDGE, MASSACHUSETTS 02139

098

## REPORT DOCUMENTATION PAGE

1a. REPORT SECURITY CLASSIFICATION Unclassified			1b. RESTRICTIVE MARKINGS		
2a. SECURITY CLASSIFICATION AUTHORITY			3. DISTRIBUTION/AVAILABILITY OF REPORT Approved for public release; distribution is unlimited.		
2b. DECLASSIFICATION/DOWNGRADING SCHEDULE					
4. PERFORMING ORGANIZATION REPORT NUMBER(S) MIT/LCS/TM-406			5. MONITORING ORGANIZATION REPORT NUMBER(S) N00014-83-K-0125 and N00014-85-K-0168		
6a. NAME OF PERFORMING ORGANIZATION MIT Laboratory for Computer Science		6b. OFFICE SYMBOL (If applicable)		7a. NAME OF MONITORING ORGANIZATION Office of Naval Research/Department of Navy	
6c. ADDRESS (City, State, and ZIP Code) 545 Technology Square Cambridge, MA 02139				7b. ADDRESS (City, State, and ZIP Code) Information Systems Program Arlington, VA 22217	
8a. NAME OF FUNDING/SPONSORING ORGANIZATION DARPA/DOD		8b. OFFICE SYMBOL (If applicable)		9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER	
8c. ADDRESS (City, State, and ZIP Code) 1400 Wilson Boulevard Arlington, VA 22217		10. SOURCE OF FUNDING NUMBERS			
		PROGRAM ELEMENT NO.		PROJECT NO.	TASK NO.
					WORK UNIT ACCESSION NO.
11. TITLE (Include Security Classification) <u>On the Correctness of Orphan Management Algorithms</u>					
12. PERSONAL AUTHOR(S) Herlihy, M., Lynch, N., Merritt, M., and Weihl, W.					
13a. TYPE OF REPORT Technical		13b. TIME COVERED FROM _____ TO _____		14. DATE OF REPORT (Year, Month, Day) 1989 August	
15. PAGE COUNT 52					
16. SUPPLEMENTARY NOTATION					
17. COSATI CODES			18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) orphans, transactions, atomic transactions, databases, distributed computing, distributed databases, serializability, concurrency control, fault-tolerance		
FIELD	GROUP	SUB-GROUP			
19. ABSTRACT (Continue on reverse if necessary and identify by block number) In a distributed system, node failures, network delays and other unpredictable occurrences can result in orphan computations--subcomputations that continue to run but whose results are no longer needed. Several algorithms have been proposed to prevent such computations from seeing inconsistent states of the shared data. In this paper, two such orphan management algorithms are analyzed. The first is an algorithm proposed at Carnegie-Mellon. The algorithms are described formally, and complete proofs of their correctness are given.  The proofs show that the fundamental concepts underlying the two algorithms are very similar in that each can be regarded as an implementation of the same high-level algorithm. By exploiting properties of information flow within transaction management systems, the algorithms ensure that orphans only see states of the shared data that they could also see if they were not orphans. When the algorithms are used in combination with any correct concurrency control algorithm, they guarantee that all computations, orphan as well as non-orphan, see consistent states of the shared data.					
20. DISTRIBUTION/AVAILABILITY OF ABSTRACT <input checked="" type="checkbox"/> UNCLASSIFIED/UNLIMITED <input type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS			21. ABSTRACT SECURITY CLASSIFICATION Unclassified		
22a. NAME OF RESPONSIBLE INDIVIDUAL Judy Little, Publications Coordinator			22b. TELEPHONE (Include Area Code) (617) 253-5894		22c. OFFICE SYMBOL

# On the Correctness of Orphan Management Algorithms

Maurice Herlihy<sup>1</sup>  
Nancy Lynch<sup>2</sup>  
Michael Merritt<sup>3</sup>  
William Weihl<sup>4</sup>

## Abstract

In a distributed system, node failures, network delays and other unpredictable occurrences can result in *orphan* computations—subcomputations that continue to run but whose results are no longer needed. Several algorithms have been proposed to prevent such computations from seeing inconsistent states of the shared data. In this paper, two such orphan management algorithms are analyzed. The first is an algorithm implemented in the Argus distributed computing system at MIT, and the second is an algorithm proposed at Carnegie-Mellon. The algorithms are described formally, and complete proofs of their correctness are given.

The proofs show that the fundamental concepts underlying the two algorithms are very similar in that each can be regarded as an implementation of the same high-level algorithm. By exploiting properties of information flow within transaction management systems, the algorithms ensure that orphans only see states of the shared data that they could also see if they were not orphans. When the algorithms are used in combination with any correct concurrency control algorithm, they guarantee that all computations, orphan as well as non-orphan, see consistent states of the shared data.

**Keywords:** orphans, transactions, atomic transactions, databases, distributed computing, distributed databases, serializability, concurrency control, fault-tolerance.

<sup>1</sup>Sponsored by the Defense Advanced Research Projects Agency (DOD), ARPA Order No. 4976 (Amendment 20), under Contracts F33615-84-K-1520 and F33615-87-C-1499, monitored by the Avionics Laboratory, Air Force Wright Aeronautical Laboratories, Wright-Patterson AFB. Address: CMU Department of Computer Science, Pittsburgh, PA.

<sup>2</sup>Supported by the National Science Foundation under Grants DCR-83-02391 and CCR-86-11442, the Defense Advanced Research Projects Agency (DARPA) under Contract N00014-83-K-0125, the Office of Naval Research under Contract N00014-85-K-0168, and the Office of Army Research under Contract DAAG29-84-K-0058. Address: MIT Laboratory for Computer Science, 545 Technology Sq., Cambridge, MA.

<sup>3</sup>Address: AT&T Bell Laboratories, Murray Hill, NJ.

<sup>4</sup>Supported by an IBM Faculty Development Award, the National Science Foundation under grants DCR-85-100014 and CCR-8716884, and the Defense Advanced Research Projects Agency (DARPA) under Contract N00014-83-K-0125. Address: MIT Laboratory for Computer Science, 545 Technology Sq., Cambridge, MA.

For	<input checked="" type="checkbox"/>
V&I	<input checked="" type="checkbox"/>
ed	<input type="checkbox"/>
tion	<input type="checkbox"/>

By _____	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	



## 1. Introduction

Nested transaction systems have been explored in a number of recent research projects (e.g., see [8, 23, 21, 1, 7]) as a means for organizing computations in distributed systems. Like ordinary transactions, nested transactions provide a simple construct for masking the effects of concurrency and failures. Nested transactions extend the usual notion of transactions [3] to permit concurrency within a single transaction. They also provide a greater degree of fault-tolerance by isolating a transaction from the failures of its descendants.

In distributed systems, various factors, including node crashes and network delays, can result in *orphan* computations—subcomputations that continue to run even though their results are no longer needed. For example, in the Argus system [8], a node making a remote request may give up because a network partition or some other problem prevents it from communicating with the other node. This may leave a process running at the called node; this process is an orphan. The orphan runs as a descendant of the transaction that made the call. Since the caller gives up by aborting the transaction that made the call, the orphan will not have any permanent effects on the observed state of the shared data.

As discussed in [9, 16, 17], even if a system is designed to prevent orphans from permanently affecting shared data, orphans are still undesirable, for two reasons. First, they waste resources: they use processor cycles, and may also hold locks, causing other computations to be delayed. Second, they may see inconsistent states of the shared data. For example, a transaction might read data at two nodes, with some invariant relating the values of the different data objects. If the transaction reads data at one of the nodes and then becomes an orphan, another transaction could change the data at both nodes before the orphan reads the data at the second node. This could happen, for example, because the first node learns that the transaction has aborted and releases its locks. While the inconsistencies seen by an orphan should not have any permanent effect on the shared data in the system, they can cause strange behavior if the orphan interacts with the external world; this can make programs difficult to design and debug.

Several algorithms have been proposed to prevent orphans from seeing inconsistent information. Early work in the area includes [19], which describes algorithms for detecting and eliminating orphans that arise because of node crashes. Nelson's work did not assume an underlying transaction mechanism, so it was difficult to assign simple semantics to abandoned computations. Recent work [24, 9, 16, 17] has studied orphans in the context of a nested transaction system, in which an abandoned computation can be *aborted*, preventing it from having any effect on the state of the system. The goal of the algorithms in [24, 9, 16, 17] is to detect and eliminate orphans before they can see inconsistent information.

### 1.1. New Results

In this paper we give formal descriptions and correctness proofs for the two orphan management algorithms in [9] and [16, 17]. The algorithm in [9] is currently in use in the Argus system.<sup>5</sup> Our proofs are completely rigorous, yet quite simple. In addition, both the

---

<sup>5</sup>Our analysis covers only orphans resulting from aborts of transactions that leave running descendants; there is another component of the Argus algorithm that handles orphans that result from node crashes in which the contents of volatile memory are destroyed

presentations and the proofs follow the intuitions that the designers have used in describing the algorithms. Although the two algorithms appear to be quite different, our proofs show that the fundamental concepts underlying them are very similar; in fact, each can be regarded as an implementation of the same high-level algorithm.

Our results relate the behavior of a system,  $S'$ , containing an orphan management algorithm to that of a corresponding system,  $S$ , having no orphan management; namely,  $S'$  must "simulate"  $S$  in the sense that each transaction in  $S'$  must see a view of the system that it could see in an execution of  $S$  in which it is not an orphan. (A transaction's "view" of the system is its sequence of interactions with the system, including the results of operations and subtransactions invoked by the transaction.) When system  $S$  includes a concurrency control algorithm that ensures that non-orphans see consistent views, our results imply that in  $S'$ , *all* transactions, orphan as well as non-orphan, see consistent views. This result provides formal justification for an informal claim sometimes made by the algorithms' designers, that the algorithms work in combination with any concurrency control algorithm.

The formal model used in this paper is based on that in [11, 12, 4]. In [11, 12], Lynch and Merritt develop a model for nested transaction systems including aborts, and use the model to show that an exclusive locking variation of Moss's algorithm [18] ensures correctness for non-orphans. The paper [4] contains improvements to the basic model in [11, 12], plus proofs that Moss' read-write algorithm and a more general commutativity-based locking algorithm also ensure correctness for non-orphans. In this paper we use the same model to describe the two orphan management algorithms mentioned above, to state correctness properties, and to prove the algorithms correct.

## 1.2. Related Work

Earlier work on verifying the Argus orphan management algorithm appears in [6]. This work is based on an earlier model for nested transaction systems that is described in [10]. The results in [6] are less general than those presented here, since they apply only to the specific concurrency control algorithm (nested locking) used by Argus. Moreover, the presentation there is much more complex than the one in this paper. Much of the complexity in [6] arises because the treatments of concurrency control and orphan management are intermingled, whereas here we are able to separate the two. The model in [11, 12, 4] provides a convenient set of concepts for describing this separation.

Other work using the model of [11, 12, 4] includes [5, 2, 20]; these papers prove correctness of algorithms for replica management, timestamp-based concurrency control and distributed transaction commit, respectively. The fact that it is possible to use the model to explain such a variety of transaction-processing algorithms is strong evidence that it is a very useful tool for modeling and analyzing nested transaction systems.

## 1.3. Organization of this Paper

The remainder of the paper is organized as follows. Section 2 contains some preliminary mathematical definitions and a brief description of *I/O automata*, which serve as the formal foundation for our work. This section may be skipped on first or cursory reading, and is included in order to make the technical presentation of this paper entirely self-contained. Section

3 contains a definition of *basic systems*, a general class of transaction-processing systems to which our results apply. These are nested transaction systems in which orphans may occur, and for which the problem of managing orphans can be precisely and intuitively stated. A basic system models the components of a nested transaction system as I/O automata. Each user program is modeled as a transaction automaton, and the rest of the system (which may include a division into objects, and may include concurrency control and recovery algorithms) is modeled as a single basic database automaton. A basic system is said to manage orphans correctly if it ensures a property called "serial correctness" for all transactions, orphans and non-orphans.

Section 4 contains some definitions and results about the dependencies among different events in a basic system; these concepts underlie the results in the rest of the paper.

Sections 5 through 8 contain the principal contributions of this paper, in which we prove the correctness of the two orphan management algorithms in [9] and [16, 17]. Our proofs have an interesting structure. We first define a simple abstract algorithm that uses global information about the history of the system, and show that it ensures that orphans see consistent views. We then formalize the Argus algorithm and the clocked algorithm from [16] in a way that only requires the use of local information, and show that each simulates the more abstract algorithm. The simulation proofs are quite simple, and do not require re-proving the properties already proved for the abstract algorithm. The correctness of the Argus and clocked algorithms then follows directly from the correctness of the abstract algorithm.

Each orphan management algorithm is described as a system obtained via simple transformations of an arbitrary basic system without orphan management. Each of these systems contains the same transactions as the given basic system, but each manages orphans using a different basic database. The abstract algorithm is modeled by the *filtered database*, which maintains information about the global history of the system, and uses tests based on this history information to prevent orphans from learning that they are orphans. The *Argus database* models the behavior of the Argus orphan management algorithm [9]; it manages orphans using tests based on local information about direct dependencies among system events. The *strictly filtered database* models another abstract algorithm, introduced to simplify the proof of the correctness of the algorithm in [16]; it also uses tests based on global history information, and is even more restrictive than the filtered database. Finally, the *clock database* models the orphan management algorithm from [16]; it manages orphans using information about logical clocks. Each of these four databases is described as the result of a transformation of the basic database.

We prove that the *filtered system* (the system consisting of the transactions and the filtered database) "simulates" the basic system in the sense that all transactions, including orphans, see a "view" that they could see in the basic system, in an execution in which they are not orphans. It follows that if the basic system ensures serial correctness for non-orphan transactions, then the filtered system ensures serial correctness for all transactions. We also prove that the *Argus system* "implements" the filtered system, and so inherits the same correctness property. Similarly, we prove that the *clock system* implements the *strictly filtered system*, which in turn implements the filtered system, thus showing that the clock system has the same correctness property as the filtered system.

Section 9 makes some of the preceding general concepts more concrete by describing two particular types of basic systems, taken from other work using this model. The first kind of basic

system, a "generic system," is appropriate for describing locking algorithms, while the second kind of basic system, a "pseudotime system," is appropriate for describing timestamp-based algorithms. Both kinds of systems specialize the notion of a basic system by splitting the basic database automaton into two kinds of components: an "object automaton" for each object in the system and a "controller automaton" that links the transactions and objects together. The concurrency control and recovery performed by the system is encapsulated within the object automata. The two kinds of systems differ in that they have slightly different interfaces between the objects and the controller. Particular information flow dependencies are described for both of these kinds of basic systems.

Section 10 contains a summary of our results and some suggestions for further work.

## 2. Formal Preliminaries

An *irreflexive partial order* is a binary relation that is irreflexive, antisymmetric and transitive.

The formal subject matter of this paper is concerned with finite and infinite sequences describing the executions of automata. Usually, we will be discussing sequences of elements from a universal set of *actions*. Formally, a *sequence*  $\beta$  of actions is a mapping from a prefix of the positive integers to the set of actions. We describe the sequence by listing the images of successive integers under the mapping, writing  $\beta = \pi_1\pi_2\pi_3\dots$ <sup>6</sup> Since the same action may occur several times in a sequence, it is convenient to distinguish the different occurrences. Thus, we refer to a particular occurrence of an action in a sequence as an *event*. Formally, an *event* in a sequence  $\beta = \pi_1\pi_2\dots$  of actions is an ordered pair  $(i, \pi)$ , where  $i$  is a positive integer and  $\pi$  is an action, such that  $\pi_i$ , the  $i^{\text{th}}$  action in  $\beta$ , is  $\pi$ .

A set of sequences  $P$  is *prefix-closed* provided that whenever  $\beta \in P$  and  $\gamma$  is a prefix of  $\beta$ , it is also the case that  $\gamma \in P$ . Similarly, a set of sequences  $P$  is *limit-closed* provided that any sequence all of whose finite prefixes are in  $P$  is also in  $P$ . We refer to any nonempty, prefix-closed and limit-closed set of sequences as a *safety property*.

### 2.1. The Input/Output Automaton Model

In order to reason carefully about complex concurrent systems such as those that implement atomic transactions, it is important to have a simple and clearly defined formal model for concurrent computation. The model we use for our work is the *input/output automaton* model [14, 15]. This model allows careful and readable descriptions of concurrent algorithms and of the correctness conditions that they are supposed to satisfy. The model can serve as the basis for rigorous proofs that particular algorithms satisfy particular correctness conditions.

This subsection contains an introduction to a simple special case of the model that is sufficient for use in this paper. In particular, in this paper we consider properties of finite executions only, and do not consider "liveness" or "fairness" properties.

Each system component is modeled as an "I/O automaton," which is a mathematical object

---

<sup>6</sup>We use the symbols  $\beta, \gamma, \dots$  for sequences of actions and the symbols  $\pi, \phi$  and  $\psi$  for individual actions.

somewhat like a traditional finite-state automaton. However, an I/O automaton need not be finite-state, but can have an infinite state set. The actions of an I/O automaton are classified as either "input," "output" or "internal." This classification is a reflection of a distinction between events (such as the receipt of a message) that are caused by the environment, events (such as sending a message) that the component can perform when it chooses and that affect the environment, and events (such as changing the value of a local variable) that a component can perform when it chooses, but that are undetectable by the environment except through their effects on later events. In the model, an automaton generates output and internal actions autonomously, and transmits output actions instantaneously to its environment. In contrast, the automaton's input is generated by the environment and transmitted instantaneously to the automaton. The distinction between input and other actions is fundamental, based on who determines when the action is performed: an automaton can establish restrictions on when it will perform an output or internal action, but it is unable to block the performance of an input action.

### 2.1.1. Action Signatures

The formal description of an automaton's actions and their classification into inputs, outputs and internal actions is given by its "action signature." An *action signature*  $S$  is an ordered triple consisting of three pairwise-disjoint sets of actions. We write  $in(S)$ ,  $out(S)$  and  $int(S)$  for the three components of  $S$ , and refer to the actions in the three sets as the *input actions*, *output actions* and *internal actions* of  $S$ , respectively. We let  $ext(S) = in(S) \cup out(S)$  and refer to the actions in  $ext(S)$  as the *external actions* of  $S$ . Also, we let  $local(S) = int(S) \cup out(S)$ , and refer to the actions in  $local(S)$  as the *locally controlled actions* of  $S$ . Finally, we let  $acts(S) = in(S) \cup out(S) \cup int(S)$ , and refer to the actions in  $acts(S)$  as the *actions* of  $S$ .

An *external action signature* is an action signature consisting entirely of external actions, that is, having no internal actions. If  $S$  is an action signature, then the *external action signature* of  $S$  is the action signature  $extsig(S) = (in(S), out(S), \emptyset)$ , i.e., the action signature that is obtained from  $S$  by removing the internal actions.

### 2.1.2. Input/Output Automata

An *input/output automaton*  $A$  (also called an *I/O automaton* or simply an *automaton*) consists of four components:<sup>7</sup>

- an action signature  $sig(A)$ ,
- a set  $states(A)$  of *states*,
- a nonempty set  $start(A) \subseteq states(A)$  of *start states*, and
- a transition relation  $steps(A) \subseteq states(A) \times acts(sig(A)) \times states(A)$ , with the property that for every state  $s'$  and input action  $\pi$  there is a transition  $(s', \pi, s)$  in  $steps(A)$ .

Note that the set of states need not be finite. We refer to an element  $(s', \pi, s)$  of  $steps(A)$  as a *step* of  $A$ . The step  $(s', \pi, s)$  is called an *input step* of  $A$  if  $\pi$  is an input action, and *output steps*, *internal steps*, *external steps* and *locally controlled steps* are defined analogously. If  $(s', \pi, s)$  is a

---

<sup>7</sup>I/O automata, as defined in [14], also include a fifth component, which is used for describing fair executions. We omit it here as it is not needed for the results described in this paper.



step of  $A$ , then  $\pi$  is said to be *enabled* in  $s'$ . Since every input action is enabled in every state, automata are said to be *input-enabled*. The input-enabling property means that an automaton is not able to block input actions. If  $A$  is an automaton, we sometimes write  $acts(A)$  as shorthand for  $acts(sig(A))$ , and likewise for  $in(A)$ ,  $out(A)$ , etc. An I/O automaton  $A$  is said to be *closed* if all its actions are locally controlled, i.e., if  $in(A) = \emptyset$ .

Note that an I/O automaton can be "nondeterministic," by which we mean two things: that more than one locally controlled action can be enabled in the same state, and that the same action, applied in the same state, can lead to different successor states. This nondeterminism is an important part of the model's descriptive power. Describing algorithms as nondeterministically as possible tends to make results about the algorithms quite general, since many results about nondeterministic algorithms apply *a fortiori* to all algorithms obtained by restricting the nondeterministic choices. Moreover, the use of nondeterminism helps to avoid cluttering algorithm descriptions and proofs with inessential details. Finally, the uncertainties introduced by asynchrony make nondeterminism an intrinsic property of real concurrent systems, and so an important property to capture in our formal model of such systems.

### 2.1.3. Executions, Schedules and Behaviors

When a system is modeled by an I/O automaton, each possible run of the system is modeled by an "execution," an alternating sequence of states and actions. The possible activity of the system is captured by the set of all possible executions that can be generated by the automaton. However, not all the information contained in an execution is important to a user of the system, nor to an environment in which the system is placed. We believe that what is important about the activity of a system is the externally visible events, and not the states or internal events. Thus, we focus on the automaton's "behaviors" — the subsequences of its executions consisting of external (i.e., input and output) actions. We regard a system as suitable for a purpose if any possible sequence of externally visible events has appropriate characteristics. Thus, in the model, we formulate correctness conditions for an I/O automaton in terms of properties of the automaton's behaviors.

Formally, an *execution fragment* of  $A$  is a finite sequence  $s_0\pi_1s_1\pi_2\dots\pi_ns_n$  or infinite sequence  $s_0\pi_1s_1\pi_2\dots\pi_ns_n\dots$  of alternating states and actions of  $A$  such that  $(s_i, \pi_{i+1}, s_{i+1})$  is a step of  $A$  for every  $i$  for which  $i+1$  exists. An execution fragment beginning with a start state is called an *execution*. We denote the set of executions of  $A$  by  $execs(A)$ , and the set of finite executions of  $A$  by  $finexecs(A)$ . A state is said to be *reachable* in  $A$  if it is the final state of a finite execution of  $A$ .

The *schedule* of an execution fragment  $\alpha$  of  $A$  is the subsequence of  $\alpha$  consisting of actions, and is denoted by  $sched(\alpha)$ . We say that  $\beta$  is a *schedule* of  $A$  if  $\beta$  is the schedule of an execution of  $A$ . We denote the set of schedules of  $A$  by  $scheds(A)$  and the set of finite schedules of  $A$  by  $finscheds(A)$ . The *behavior* of a sequence  $\beta$  of actions in  $acts(A)$ , denoted by  $beh(\beta)$ , is the subsequence of  $\beta$  consisting of actions in  $ext(A)$ . The *behavior* of an execution fragment  $\alpha$  of  $A$ , denoted by  $beh(\alpha)$ , is defined to be  $beh(sched(\alpha))$ . We say that  $\beta$  is a *behavior* of  $A$  if  $\beta$  is the behavior of an execution of  $A$ . We denote the set of behaviors of  $A$  by  $behs(A)$  and the set of finite behaviors of  $A$  by  $finbehs(A)$ .

We say that a finite schedule  $\beta$  of  $A$  *can leave*  $A$  in state  $s$  if there is some finite execution  $\alpha$  of  $A$  with final state  $s$  and with  $sched(\alpha) = \beta$ . Similarly, a finite behavior  $\beta$  of  $A$  *can leave*  $A$  in

state  $s$  if there is some finite execution  $\alpha$  of  $A$  with final state  $s$  and with  $\text{beh}(\alpha) = \beta$ . We say that an action  $\pi$  is *enabled after* a finite schedule or behavior  $\beta$  of  $A$  if there is a state  $s$  such that  $\beta$  can leave  $A$  in state  $s$  and  $\pi$  is enabled in  $s$ .

An *extended step* of an automaton  $A$  is a triple of the form  $(s', \beta, s)$ , where  $s'$  and  $s$  are in  $\text{states}(A)$ ,  $\beta$  is a finite sequence of actions in  $\text{acts}(A)$ , and there is an execution fragment of  $A$  having  $s'$  as its first state,  $s$  as its last state and  $\beta$  as its schedule.

If  $\beta$  is any sequence of actions and  $\Phi$  is a set of actions, we write  $\beta|\Phi$  to denote the subsequence of  $\beta$  containing all occurrences of actions in  $\Phi$ . If  $A$  is an automaton, we write  $\beta|A$  for  $\beta|\text{acts}(A)$ .

## 2.2. Composition

Often, a single system can also be viewed as a combination of several component systems interacting with one another. To reflect this in our model, we define a "composition" operation by which several I/O automata can be combined to yield a single I/O automaton. Our composition operator connects each output action of the component automata with the identically named input actions of any number (usually one) of the other component automata. In the resulting system, an output action is generated autonomously by one component and is thought of as being instantaneously transmitted to all components having the same action as an input. All such components are passive recipients of the input, and take steps simultaneously with the output step.

### 2.2.1. Composition of Action Signatures

We first define composition of action signatures. Let  $I$  be an index set that is at most countable. A collection  $\{S_i\}_{i \in I}$  of action signatures is said to be *strongly compatible*<sup>8</sup> if the following properties hold:

1.  $\text{out}(S_i) \cap \text{out}(S_j) = \emptyset$  for all  $i, j \in I$  such that  $i \neq j$ ,
2.  $\text{int}(S_i) \cap \text{acts}(S_j) = \emptyset$  for all  $i, j \in I$  such that  $i \neq j$ , and
3. no action is in  $\text{acts}(S_i)$  for infinitely many  $i$ .

Thus, no action is an output of more than one signature in the collection, and internal actions of any signature do not appear in any other signature in the collection. Moreover, we do not permit actions involving infinitely many component signatures.

The *composition*  $S = \prod_{i \in I} S_i$  of a collection of strongly compatible action signatures  $\{S_i\}_{i \in I}$  is defined to be the action signature with

- $\text{in}(S) = \cup_{i \in I} \text{in}(S_i) - \cup_{i \in I} \text{out}(S_i)$ ,
- $\text{out}(S) = \cup_{i \in I} \text{out}(S_i)$ , and
- $\text{int}(S) = \cup_{i \in I} \text{int}(S_i)$ .

Thus, output actions are those that are outputs of any of the component signatures, and similarly

---

<sup>8</sup>A weaker notion called "compatibility" is defined in [14], consisting of the first two of the three given properties only. For the purposes of this paper, only the stronger notion will be required.

for internal actions. Input actions are any actions that are inputs to any of the component signatures, but outputs of no component signature.

### 2.2.2. Composition of Automata

A collection  $\{A_i\}_{i \in I}$  of automata is said to be *strongly compatible* if their action signatures are strongly compatible. The *composition*  $A = \prod_{i \in I} A_i$  of a strongly compatible collection of automata  $\{A_i\}_{i \in I}$  has the following components:<sup>9</sup>

- $\text{sig}(A) = \prod_{i \in I} \text{sig}(A_i)$ ,
- $\text{states}(A) = \prod_{i \in I} \text{states}(A_i)$ ,
- $\text{start}(A) = \prod_{i \in I} \text{start}(A_i)$ , and
- $\text{steps}(A)$  is the set of triples  $(s', \pi, s)$  such that for all  $i \in I$ , (a) if  $\pi \in \text{acts}(A_i)$  then  $(s'[i], \pi, s[i]) \in \text{steps}(A_i)$ , and (b) if  $\pi \notin \text{acts}(A_i)$  then  $s'[i] = s[i]$ .<sup>10</sup>

Since the automata  $A_i$  are input-enabled, so is their composition, and hence their composition is an automaton. Each step of the composition automaton involves all the automata that have a particular action in their action signature performing that action concurrently, while the automata that do not have that action in their signature do nothing. We will often refer to an automaton formed by composition as a "system" of automata.

If  $\alpha = s_0 \pi_1 s_1 \dots$  is an execution of  $A$ , let  $\alpha|A_i$  be the sequence obtained by deleting  $\pi_j s_j$  when  $\pi_j$  is not an action of  $A_i$ , and replacing the remaining  $s_j$  by  $s_j[i]$ . Recall that we have previously defined a projection operator for action sequences. The two projection operators are related in the obvious way:  $\text{sched}(\alpha|A_i) = \text{sched}(\alpha)|A_i$ , and similarly  $\text{beh}(\alpha|A_i) = \text{beh}(\alpha)|A_i$ .

In the course of our discussions we will often reason about automata without specifying their internal actions. To avoid tedious arguments about compatibility, henceforth we assume that unspecified internal actions of any automaton are unique to that automaton, and do not occur as internal or external actions of any of the other automata we discuss.

All of the systems that we will use for modeling transactions are closed systems, that is, each action is an output of some component. Also, each output of a component will be an input of at most one other component.

### 2.2.3. Properties of Systems of Automata

Here we give basic results relating executions, schedules and behaviors of a system of automata to those of the automata being composed. The first result says that the projections of executions of a system onto the components are executions of the components, and similarly for schedules, etc.

**Proposition 1:** Let  $\{A_i\}_{i \in I}$  be a strongly compatible collection of automata, and let

---

<sup>9</sup>Note that the second and third components listed are just ordinary Cartesian products, while the first component uses the previous definition of composition of action signatures.

<sup>10</sup>We use the notation  $s[i]$  to denote the  $i^{\text{th}}$  component of the state vector  $s$ .

$A = \prod_{i \in I} A_i$ . If  $\alpha \in \text{execs}(A)$  then  $\alpha|A_i \in \text{execs}(A_i)$  for all  $i \in I$ . Moreover, the same result holds for  $\text{finexecs}$ ,  $\text{scheds}$ ,  $\text{finscheds}$ ,  $\text{behs}$  and  $\text{finbehs}$  in place of  $\text{execs}$ .

Converses can also be proved for all the parts of the preceding proposition. The following are most useful. They say that schedules and behaviors of component automata can be "patched together" to form schedules and behaviors of the composition.

**Proposition 2:** Let  $\{A_i\}_{i \in I}$  be a strongly compatible collection of automata, and let  $A = \prod_{i \in I} A_i$ .

1. Let  $\beta$  be a sequence of actions in  $\text{acts}(A)$ . If  $\beta|A_i \in \text{scheds}(A_i)$  for all  $i \in I$ , then  $\beta \in \text{scheds}(A)$ .
2. Let  $\beta$  be a finite sequence of actions in  $\text{acts}(A)$ . If  $\beta|A_i \in \text{finscheds}(A_i)$  for all  $i \in I$ , then  $\beta \in \text{finscheds}(A)$ .
3. Let  $\beta$  be a sequence of actions in  $\text{ext}(A)$ . If  $\beta|A_i \in \text{behs}(A_i)$  for all  $i \in I$ , then  $\beta \in \text{behs}(A)$ .
4. Let  $\beta$  be a finite sequence of actions in  $\text{ext}(A)$ . If  $\beta|A_i \in \text{finbehs}(A_i)$  for all  $i \in I$ , then  $\beta \in \text{finbehs}(A)$ .

The preceding proposition is useful in proving that a sequence of actions is a behavior of a system  $A$ : it suffices to show that the sequence's projections are behaviors of the components of  $A$  and then to appeal to Proposition 2.

### 2.3. Implementation

We define a notion of "implementation" of one automaton by another. Let  $A$  and  $B$  be automata with the same external action signature, i.e., with  $\text{extsig}(A) = \text{extsig}(B)$ . Then  $A$  is said to *implement*  $B$  if  $\text{finbehs}(A) \subseteq \text{finbehs}(B)$ . One way in which this notion can be used is the following. Suppose we can show that an automaton  $B$  is "correct," in the sense that its finite behaviors all satisfy some specified property. Then if another automaton  $A$  implements  $B$ ,  $A$  is also correct. One can also show that if  $A$  implements  $B$ , then replacing  $B$  by  $A$  in any system yields a new system in which all finite behaviors are behaviors of the original system.

One useful technique for showing that one automaton implements another is to give a correspondence between states of the two automata. Such a correspondence can often be expressed in the form of a kind of abstraction mapping that we call a "possibilities mapping," defined as follows. Suppose  $A$  and  $B$  are automata with the same external action signature, and suppose  $f$  is a mapping from  $\text{states}(A)$  to the power set of  $\text{states}(B)$ . That is, if  $s$  is a state of  $A$ ,  $f(s)$  is a set of states of  $B$ . The mapping  $f$  is said to be a *possibilities mapping* from  $A$  to  $B$  if the following conditions hold:

1. For every start state  $s_0$  of  $A$ , there is a start state  $t_0$  of  $B$  such that  $t_0 \in f(s_0)$ .
2. Let  $s'$  be a reachable state of  $A$ ,  $t' \in f(s')$  a reachable state of  $B$ , and  $(s', \pi, s)$  a step of  $A$ . Then there is an extended step  $(t', \gamma, t)$  of  $B$  (possibly having an empty schedule) such that the following conditions are satisfied:
  - a.  $\gamma \text{ext}(B) = \pi \text{ext}(A)$ , and
  - b.  $t \in f(s)$ .

The following proposition shows that giving a possibilities mapping from A to B is sufficient to show that A implements B.

**Proposition 3:** Suppose that A and B are automata with the same external action signature and there is a possibilities mapping,  $f$ , from A to B. Then A implements B.

## 2.4. Preserving Properties

Although automata in our model are unable to block input actions, it is often convenient to restrict attention to those behaviors in which the environment provides inputs in a "sensible" way, that is, where the environment obeys certain "well-formedness" restrictions. A useful way of discussing such restrictions is in terms of the notion that an automaton "preserves" a property of behaviors: as long as the environment does not violate the property, neither does the automaton. Such a notion is primarily interesting for safety properties. Let  $\Phi$  be a set of actions and  $P$  a safety property for sequences of actions in  $\Phi$ . Let  $A$  be an automaton with  $\Phi \cap \text{int}(A) = \emptyset$ . We say that  $A$  *preserves*  $P$  if  $\beta\pi|\Phi \in P$  whenever  $\beta|\Phi \in P$ ,  $\pi \in \text{out}(A)$  and  $\beta\pi|A \in \text{finbehs}(A)$ .

Thus, if an automaton preserves a property  $P$ , the automaton is not the first to violate  $P$ : as long as the environment only provides inputs such that the cumulative behavior satisfies  $P$ , the automaton will only perform outputs such that the cumulative behavior satisfies  $P$ . In many cases of interest, we will have  $\Phi \subseteq \text{ext}(A)$ ; note that even in this case, the fact that an automaton  $A$  preserves  $P$  does not imply that all of  $A$ 's behaviors, when restricted to  $\Phi$ , satisfy  $P$ . It is possible for a behavior of  $A$  to fail to satisfy  $P$ , if an input causes a violation of  $P$ . However, the following proposition gives a way to deduce that all of a system's behaviors satisfy  $P$ . The proposition says that, under certain conditions, if all components of a system preserve  $P$ , then all the behaviors of the composition satisfy  $P$ .

**Proposition 4:** Let  $\{A_i\}_{i \in I}$  be a strongly compatible collection of automata and let  $A = \prod_{i \in I} A_i$ . Let  $\Phi$  be a set of actions such that  $\Phi \cap \text{int}(A) = \emptyset$ , and let  $P$  be a safety property for actions in  $\Phi$ . Suppose that for each  $i \in I$ ,  $A_i$  preserves  $P$ . Then  $A$  preserves  $P$ . Furthermore, if  $\Phi \cap \text{in}(A) = \emptyset$ , then  $\text{behs}(A)|\Phi \subseteq P$ . That is, if  $\beta \in \text{behs}(A)$ , then  $\beta|\Phi \in P$ .

**Proof:** Let  $\beta$  be a sequence of actions such that  $\beta|\Phi \in P$ ,  $\pi \in \text{out}(A)$  and  $\beta\pi|A \in \text{finbehs}(A)$ . Then  $\pi \in \text{out}(A_i)$  for some  $i \in I$ , and  $\beta\pi|A_i \in \text{finbehs}(A_i)$ , by Proposition 2. Since  $A_i$  preserves  $P$ ,  $\beta\pi|\Phi \in P$ .

Now suppose that  $\Phi \cap \text{in}(A) = \emptyset$ , and let  $\beta \in \text{behs}(A)$ . Since  $A$  preserves  $P$ , by a simple induction, every finite prefix of  $\beta|\Phi$  is in  $P$ . Then  $\beta|\Phi \in P$ , by the limit-closure of  $P$ .  $\square$

## 3. Basic Systems

In this section, we define "basic systems," the class of transaction-processing systems to which our results apply. Basic systems generalize both the generic systems of [4] and the pseudotime systems of [2]. We also define correctness conditions for basic systems, in particular, the notion of correct management of orphans.

### 3.1. Overview

Transaction-processing systems consist of user-provided transaction code, plus transaction-processing algorithms designed to coordinate the activities of different transactions. The transactions are written by application programmers in a suitable programming language. Transactions are permitted to invoke operations on data objects. In addition, if nesting is allowed, transactions can invoke subtransactions and receive responses from the subtransactions describing the results of their processing.

In a transaction-processing system, the transaction-processing algorithms interact with the transactions, making decisions about when to schedule subtransactions and operations on data objects. The transaction-processing algorithms include concurrency control and recovery algorithms. In many interesting cases (e.g., for locking algorithms), the transaction-processing algorithms can be naturally divided into a "controller" and a collection of "objects," where each object includes concurrency control and recovery algorithms appropriate for that object and the controller manages communication among the transactions and objects. We do not, however, require this division for our general results.

The transaction-processing systems studied in this paper are called "basic systems." In the organization we consider, the transaction-processing algorithms are represented by a component called a "basic database." Each component of a basic system is modeled as an I/O automaton. That is, each transaction is an automaton, and the basic database is another automaton.

The nested structure of transactions is modeled by describing each transaction and subtransaction in the transaction nesting structure as a separate I/O automaton. If a parent transaction  $T$  wishes to invoke a child transaction  $T'$ ,  $T$  issues an output action that "requests that  $T'$  be created." The basic database receives this request, and at some later time might issue an action that is an input to the child  $T'$  and corresponds to the "creation" of  $T'$ . Thus, the different transactions in the nesting structure comprise a forest of automata, communicating with each other indirectly through the basic database. The highest-level transactions, i.e., those that are not subtransactions of any other transactions, are the roots in this forest.

It is actually more convenient to model the transaction nesting structure as a tree rather than as a forest. Therefore, we add an extra "root" automaton as a "dummy transaction," located at the top of the transaction nesting structure. The highest-level user-defined transactions are considered to be children of this new root. The root can be thought of as modeling the outside world, from which invocations of top-level transactions originate and to which reports about the results of such transactions are sent.

In the rest of this section, we define "basic systems" and state the correctness conditions that they are supposed to satisfy.

### 3.2. System Types

We begin by defining a type structure that will be used to name the transactions and objects in a basic system.

A *system type* consists of the following:

- a set  $T$  of *transaction names*,

- a distinguished transaction name  $T_0 \in \mathcal{T}$ ,
- a subset *accesses* of  $\mathcal{T}$  not containing  $T_0$ ,
- a mapping *parent*:  $\mathcal{T} - \{T_0\} \rightarrow \mathcal{T}$ , which configures the set of transaction names into a tree, with  $T_0$  as the root and the accesses as the leaves,
- a set  $X$  of *object names*,
- a mapping *object*: *accesses*  $\rightarrow X$ , and
- a set  $V$  of *return values*.

Each element of the set "accesses" is called an *access* transaction name, or simply an *access*. Also, if  $\text{object}(T) = X$  we say that  $T$  is an *access to X*.

In referring to the transaction tree, we use standard tree terminology, such as "leaf node," "internal node," "child," "ancestor," and "descendant." As a special case, we consider any node to be its own ancestor and its own descendant, i.e., the "ancestor" and "descendant" relations are reflexive. We also use the notion of a "least common ancestor" of two nodes.

The transaction tree describes the nesting structure for transaction names, with  $T_0$  as the name of the dummy "root transaction." Each child node in this tree represents the name of a subtransaction of the transaction named by its parent. The children of  $T_0$  represent names of the top-level user-defined transactions. The accesses represent names for the lowest-level transactions in the transaction nesting structure; we will use these lowest-level transaction names to model operations on data objects. Thus, the only transactions that actually access data are the leaves of the transaction tree and these do nothing else. The internal nodes model transactions whose function is to create and manage subtransactions including accesses, but they do not access data directly.

The tree structure should be thought of as a predefined naming scheme for all possible transactions that might ever be invoked. In any particular execution, however, only some of these transactions will actually take steps. We imagine that the tree structure is known in advance by all components of a system. The tree will, in general, be an infinite structure with infinite branching.

The set  $X$  is the set of names for the objects used in the system. Each access transaction name is assumed to be an access to some particular object, as designated by the "object" mapping. The set  $V$  of return values is the set of possible values that might be returned by successfully completed transactions to their parent transactions.

For the rest of this paper, we will fix a particular system type.

### 3.3. General Structure of Basic Systems

A basic system for a given system type is a closed system consisting of a "transaction automaton"  $A_T$  for each non-access transaction name  $T$  and a single "basic database automaton"  $B$ . Later in this section, we will give conditions to be satisfied by the transaction and basic database automata. Here, we just describe the signatures of these automata, in order to explain how the automata are interconnected.

Figure 1 depicts the structure of a basic system.

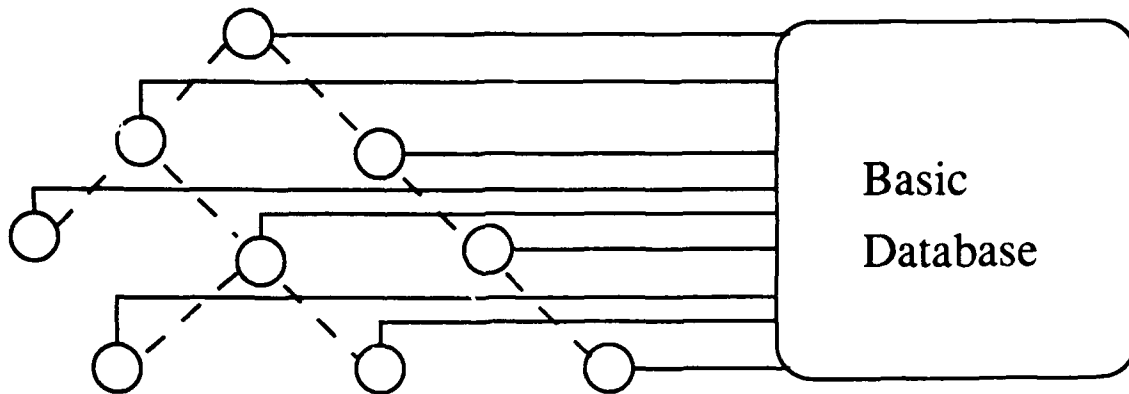


Figure 1: Basic System

The transaction nesting structure is indicated in part by dotted lines between transaction automata corresponding to parent and child. Access transactions do not have associated automata, and so the diagram does not indicate the parents of accesses. The direct connections between automata (via shared actions) are indicated by solid lines. Thus, the transaction automata interact directly with the basic database, but not directly with each other.

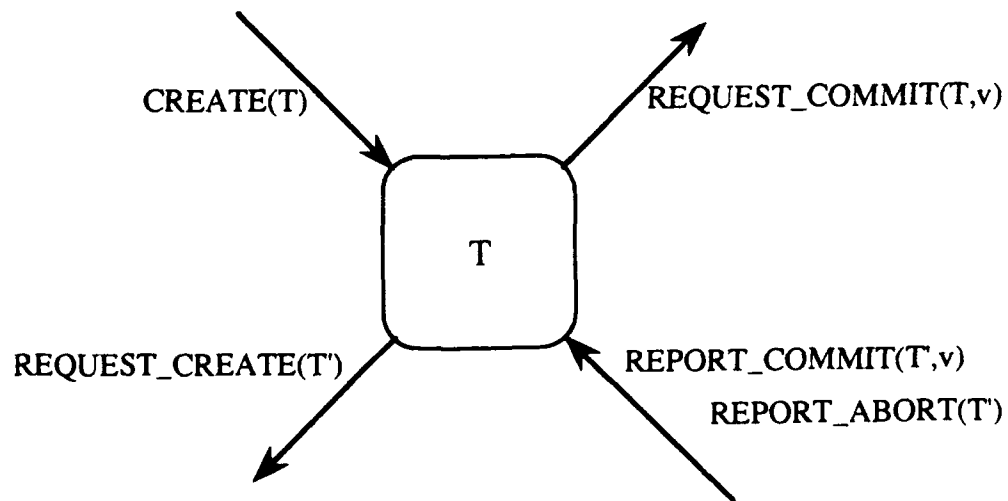


Figure 2: Transaction Interface

Figure 2 shows the interface of a transaction automaton in more detail. The automaton for transaction name  $T$  has an input action  $CREATE(T)$ , which is generated by the basic database in order to initiate  $T$ 's processing. We do not include explicit arguments to a transaction in our model; rather we suppose that there is a different transaction for each possible set of arguments,



and so any input to the transaction is encoded in the name of the transaction.  $T$  has  $REQUEST\_CREATE(T')$  actions for each child  $T'$  of  $T$  in the transaction nesting structure; these are requests for creation of child transactions, and are communicated directly to the basic database. At some later time, the basic database might respond to a  $REQUEST\_CREATE(T')$  action by issuing a  $CREATE(T')$  action; in case  $T'$  is not an access, this action is an input to the automaton for transaction  $T'$ .  $T$  also has  $REPORT\_COMMIT(T',v)$  and  $REPORT\_ABORT(T')$  input actions, by which the basic database informs  $T$  about the fate (commit or abort) of its previously requested child  $T'$ . In the case of a commit, the report includes a return value  $v$  that provides information about the activity of  $T'$ ; in the case of an abort, no information is returned. Finally,  $T$  has a  $REQUEST\_COMMIT(T,v)$  output action, by which it announces to the basic database that it has completed its activity successfully, with a particular result that is described by return value  $v$ .

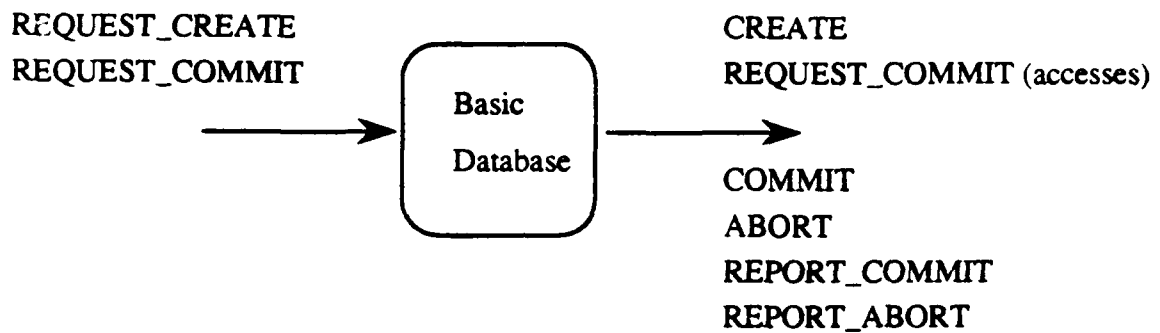


Figure 3: Basic Database Interface

Figure 3 shows the basic database interface. The basic database in any particular basic system receives the previously mentioned  $REQUEST\_CREATE$  and  $REQUEST\_COMMIT$  actions as inputs from the transaction automata. It produces  $CREATE$  actions as outputs, thereby awakening transaction automata or invoking operations on objects. The basic database also produces  $REQUEST\_COMMIT(T,v)$  output actions for accesses  $T$ ; these represent responses to the invocations of operations on objects. The value  $v$  in a  $REQUEST\_COMMIT(T,v)$  action is a return value returned by the operation as part of its response. The basic database also produces  $COMMIT(T)$  and  $ABORT(T)$  actions for arbitrary transaction names  $T \neq T_0$ , representing decisions about whether the designated transaction commits or aborts. For technical convenience, we classify the  $COMMIT$  and  $ABORT$  actions and the  $REQUEST\_COMMIT$  and  $CREATE$  actions for access transactions as output actions of the basic database, even when they are not inputs to any other system component.<sup>11</sup> The basic database also has  $REPORT\_COMMIT$  and  $REPORT\_ABORT$  actions as outputs, by which it communicates the fates of transactions to their parents.

Different basic databases may include additional output actions. The final section of this paper

<sup>11</sup>Classifying actions as outputs even though they are not inputs to any other system component is permissible in the I/O automaton model. In this case, it would also be possible to classify these actions as internal actions of the basic database, but then the statements and proofs of the ensuing results would be more complicated.

describes "generic databases," which have additional outputs by which the fates of transactions are communicated to the objects, so that locks may be released. "Pseudotime databases" are a second example, which contain additional outputs involving the management of timestamp data.

As is always the case for I/O automata, the components of a system are determined statically. Even though we referred earlier to the action of "creating" a child transaction, the model treats the child transaction as if it had been there all along. The CREATE action is treated formally as an input action to the child transaction; the child transaction will be constrained not to perform any output actions until such a CREATE action occurs. A consequence of this method of modeling dynamic creation of transactions is that the system must include automata for all possible transactions that might ever be created in any execution. In most interesting cases, this means that the system will include infinitely many transaction automata.

In our work, it is convenient to use two separate actions, REQUEST\_CREATE and CREATE, to describe what happens when a subtransaction is activated. This separation occurs in actual distributed systems such as Argus, and is important in our results and proofs. Similar remarks hold for the distinction among REQUEST\_COMMIT, COMMIT and REPORT\_COMMIT actions.

### 3.4. Serial Actions

The external actions of a basic system of a given system type include the *serial actions* for that type. The *serial actions* for a given system type are defined to be the actions listed in the preceding subsection: CREATE(T) and REQUEST\_COMMIT(T,v), where T is any transaction name and v is a return value, and REQUEST\_CREATE(T), COMMIT(T), ABORT(T), REPORT\_COMMIT(T,v), and REPORT\_ABORT(T), where  $T \neq T_0$  is a transaction name and v is a return value.<sup>12</sup> If  $\beta$  is a sequence of actions, define *serial*( $\beta$ ) to be the subsequence of  $\beta$  containing all the serial actions in  $\beta$ .

In this subsection, we define some simple concepts involving serial actions. All the definitions in this subsection are based on the set of actions only, and not on the specific automata in any particular system. For this reason, we present these definitions here, before going on (in the next subsection) to give more information about the basic system components.

We first present some fundamental definitions, and then we define notions of "well-formedness" for sequences of actions.

#### 3.4.1. Terminology

The COMMIT(T) and ABORT(T) actions are called *completion* actions for T, while the REPORT\_COMMIT(T,v) and REPORT\_ABORT(T) actions are called *report* actions for T.

We associate transaction names with some of the serial actions, as follows. Let T be a transaction name. If  $\pi$  is either a CREATE(T) or a REQUEST\_COMMIT(T,v) action, or is a REQUEST\_CREATE(T'), REPORT\_COMMIT(T',v') or REPORT\_ABORT(T'), where T' is a

---

<sup>12</sup>These actions are called "serial actions" because they are exactly the external actions of a "serial system" of the given type. More will be said about serial systems later in the paper.

child of  $T$ , then we define  $transaction(\pi)$  to be  $T$ . If  $\pi$  is a completion action, then  $transaction(\pi)$  is undefined. In some contexts, we will also need to associate a transaction with completion actions; since a completion action for  $T$  can be thought of as occurring in between  $T$  and  $parent(T)$ , some of the time we will want to associate  $T$  with the action, and at other times we will want to associate  $parent(T)$  with it. Thus, we extend the " $transaction(\pi)$ " definition in two different ways. If  $\pi$  is any serial action, then we define  $hightransaction(\pi)$  to be  $transaction(\pi)$  if  $\pi$  is not a completion action, and to be  $parent(T)$ , if  $\pi$  is a completion action for  $T$ . Also, if  $\pi$  is any serial action, we define  $lowtransaction(\pi)$  to be  $transaction(\pi)$  if  $\pi$  is not a completion action, and to be  $T$ , if  $\pi$  is a completion action for  $T$ . In particular,  $hightransaction(\pi) = lowtransaction(\pi) = transaction(\pi)$  for all serial actions  $\pi$  for which  $transaction(\pi)$  is defined.

We also require notation for the object associated with any serial action whose transaction is an access. If  $\pi$  is a serial action of the form  $CREATE(T)$  or  $REQUEST\_COMMIT(T,v)$ , where  $T$  is an access to  $X$ , then we define  $object(\pi)$  to be  $X$ .

We extend the preceding notation to events as well as actions. For example, if  $\pi$  is an event, then we write  $transaction(\pi)$  to denote the transaction of the action of which  $\pi$  is an occurrence. We extend the definitions of "hightransaction," "lowtransaction," and "object" similarly. We will extend other notation in this paper in the same way, without further explanation.

Now we require terminology to describe the status of a transaction during execution. Let  $\beta$  be a sequence of actions. A transaction name  $T$  is said to be *active* in  $\beta$  provided that  $\beta$  contains a  $CREATE(T)$  event but no  $REQUEST\_COMMIT$  event for  $T$ . Similarly,  $T$  is said to be *live* in  $\beta$  provided that  $\beta$  contains a  $CREATE(T)$  event but no completion event for  $T$ . (However, note that  $\beta$  may contain a  $REQUEST\_COMMIT$  for  $T$ .) Also,  $T$  is said to be an *orphan* in  $\beta$  if there is an  $ABORT(U)$  action in  $\beta$  for some ancestor  $U$  of  $T$ .

We have already used projection operators to restrict action sequences to particular sets of actions, and to actions of particular automata. We now introduce another projection operator, this time to sets of transaction names. Namely, if  $\beta$  is a sequence of actions and  $\mathcal{U}$  is a set of transaction names, then  $\beta|\mathcal{U}$  is defined to be the sequence  $\beta|\{\pi: transaction(\pi) \in \mathcal{U}\}$ . If  $T$  is a transaction name, we sometimes write  $\beta|T$  as shorthand for  $\beta|\{T\}$ . Similarly, if  $\beta$  is a sequence of actions and  $X$  is an object name, we sometimes write  $\beta|X$  to denote  $\beta|\{\pi: object(\pi) = X\}$ .

### 3.4.2. Well-Formedness

We will place very few constraints on the transaction automata and basic database automaton in our definition of a basic system. However, we will want to assume that certain simple properties are guaranteed; for example, a transaction should not take steps until it has been created, and the basic database should not create a transaction that has not been requested. Such requirements are captured by well-formedness conditions, which are fundamental safety properties of sequences of external actions of the transaction and basic database automata. We define those conditions here.

First, we define "transaction well-formedness." Let  $T$  be any transaction name. A sequence  $\beta$  of serial actions  $\pi$  with  $transaction(\pi) = T$  is defined to be *transaction well-formed* for  $T$  provided the following conditions hold.

1. The first event in  $\beta$ , if any, is a  $CREATE(T)$  event, and there are no other  $CREATE$  events.

2. There is at most one REQUEST\_CREATE( $T'$ ) event in  $\beta$  for each child  $T'$  of  $T$ .
3. Any report event for a child  $T'$  of  $T$  is preceded by REQUEST\_CREATE( $T'$ ) in  $\beta$ .
4. There is at most one report event in  $\beta$  for each child  $T'$  of  $T$ .
5. If a REQUEST\_COMMIT event for  $T$  occurs in  $\beta$ , then it is preceded by a report event for each child  $T'$  of  $T$  for which there is a REQUEST\_CREATE( $T'$ ) in  $\beta$ .
6. If a REQUEST\_COMMIT event for  $T$  occurs in  $\beta$ , then it is the last event in  $\beta$ .

In particular, if  $T$  is an access, then the only sequences that are transaction well-formed for  $T$  are the prefixes of the two-event sequences of the form CREATE( $T$ )REQUEST\_COMMIT( $T, v$ ). For any  $T$ , it is easy to see that the set of transaction well-formed sequences for  $T$  is a safety property, i.e., that it is prefix-closed and limit-closed.

Next, we define "basic database well-formedness." A sequence  $\beta$  of serial actions is defined to be *basic database well-formed* provided the following conditions hold.

1. The sequence  $\beta|T$  is transaction well-formed, for all transaction names  $T$ .
2. If a CREATE( $T$ ) event occurs in  $\beta$ , for  $T \neq T_0$ , then there is a preceding REQUEST\_CREATE( $T$ ) in  $\beta$ .
3. If there is a COMMIT( $T$ ) event in  $\beta$ , then there is a preceding REQUEST\_COMMIT( $T, v$ ) event in  $\beta$ , for some  $v$ .
4. If there is an ABORT( $T$ ) event in  $\beta$ , then there is a preceding REQUEST\_CREATE( $T$ ) event in  $\beta$ .
5. There is at most one completion event in  $\beta$  for each transaction name  $T$ .
6. If there is a REPORT\_COMMIT( $T, v$ ) event in  $\beta$ , then there is a preceding REQUEST\_COMMIT( $T, v$ ) event in  $\beta$  and a preceding COMMIT( $T$ ) event in  $\beta$ .
7. If there is a REPORT\_ABORT( $T$ ) event in  $\beta$ , then there is a preceding ABORT( $T$ ) event in  $\beta$ .
8. There is at most one report event in  $\beta$  for each transaction name  $T$ .

### 3.5. Basic Systems

We are now ready to define "basic systems." Basic systems are composed of transaction automata, one for each non-access transaction name, and a single basic database automaton. We describe the two kinds of components in turn.

#### 3.5.1. Transaction Automata

A *transaction automaton*  $A_T$  for a non-access transaction name  $T$  (of the given system type) is an I/O automaton with the following external action signature.

Input:

CREATE(T)  
 REPORT\_COMMIT(T',v), for every child T' of T, and return value v  
 REPORT\_ABORT(T'), for every child T' of T

Output:

REQUEST\_CREATE(T'), for every child T' of T  
 REQUEST\_COMMIT(T,v), for every return value v

In addition,  $A_T$  may have an arbitrary set of internal actions. We require  $A_T$  to preserve transaction well-formedness for T, as defined in the preceding subsection. As discussed earlier, this does not mean that all behaviors of  $A_T$  are transaction well-formed, but it does mean that as long as the environment of  $A_T$  does not violate transaction well-formedness,  $A_T$  will not do so. Except for that requirement, transaction automata can be chosen arbitrarily. Note that if  $\beta$  is a sequence of actions, then  $\beta \upharpoonright T = \beta \text{ext}(A_T)$ .

Transaction automata are intended to be general enough to model the transactions defined in any reasonable programming language. Particular programming languages may impose additional restrictions on transaction behavior. (For example, Argus suspends activity in transactions until subtransactions complete.) However, our results do not require such restrictions.

### 3.5.2. Basic Database Automata

A *basic database automaton* is also modeled as an I/O automaton. A basic database passes requests for the creation of subtransactions to the appropriate recipient, initiates REQUEST\_COMMIT actions for accesses, makes decisions about the commit or abort of transactions, and passes reports about the completion of children back to their parents. It may also carry out other activity.

A basic database has the following actions in its external action signature.

Input:

REQUEST\_CREATE(T),  $T \neq T_0$   
 REQUEST\_COMMIT(T,v), T a non-access transaction name

Output:

CREATE(T)  
 REQUEST\_COMMIT(T,v), T an access transaction name  
 COMMIT(T),  $T \neq T_0$   
 ABORT(T),  $T \neq T_0$   
 REPORT\_COMMIT(T,v),  $T \neq T_0$   
 REPORT\_ABORT(T),  $T \neq T_0$

In addition, it may have other arbitrary output actions, as well as arbitrary internal actions. Depending upon the design of the particular basic database automaton, some of the additional output actions may be associated with particular objects. Hence, each basic database is assumed to come equipped with an extension of the *object* partial mapping on actions, which may associate some of these additional, non-serial output actions with particular object names. That is, each non-serial output action  $\pi$  may (but need not) have  $\text{object}(\pi)$  defined.

The REQUEST\_CREATE and REQUEST\_COMMIT inputs are intended to be identified with the corresponding outputs of transaction automata, and conversely, all the CREATE and report outputs (except those CREATE(T) actions for which T is an access) are identified with the corresponding inputs of transaction automata. A basic database is required to preserve basic database well-formedness.

There are many examples of basic databases in the literature. For example, the composition of the generic controller and generic objects of [4] preserves basic database well-formedness, and so is an example of a basic database. The same is true for the composition of the pseudotime controller and pseudotime objects of [2]. We will present these examples in more detail later in this paper. In fact, we claim that almost all interesting transaction-processing algorithms can be modeled as basic databases. (See [13] for additional examples.)

Our notion of basic database identifies the aspects of transaction-processing algorithms that are relevant to our analysis of orphan management algorithms. It turns out that the details of how synchronization and recovery are implemented by a basic database are largely irrelevant. Indeed, this is one of the important contributions of this paper: we are able to state correctness conditions for and verify orphan management algorithms in a way that is independent of the concurrency control and recovery methods used within the basic database.

### 3.5.3. Basic Systems

A *basic system* B is the composition of a strongly compatible set of automata indexed by the union of the set of non-access transaction names and the singleton set {BD} (for "basic database"). Associated with each non-access transaction name T is a transaction automaton  $A_T$  for T. Associated with the name BD is a basic database automaton for the system type.

When the particular basic system B is understood from the context, we call its external actions the *basic actions*, and its executions, schedules and behaviors the *basic executions*, *basic schedules* and *basic behaviors*, respectively. The following proposition says that basic behaviors have the appropriate well-formedness properties.

**Proposition 5:** If  $\beta$  is a basic behavior, then the following conditions hold.

1. For every transaction name T,  $\beta|T$  is transaction well-formed for T.
2. The sequence serial( $\beta$ ) is basic database well-formed.

**Proof:** Note first that the basic database preserves basic database well-formedness, and this immediately implies that it preserves transaction well-formedness for every transaction name. Next, note that each transaction automaton preserves transaction well-formedness for the appropriate transaction name. Furthermore, it has in its signature no actions of other transactions, and so preserves transaction well-formedness for all transaction names. The first part of the proposition follows by Proposition 4.

A simple induction shows that each transaction automaton also preserves basic database well-formedness, and the second conclusion follows also from Proposition 4.

□

### 3.6. Serial Correctness

In this subsection, we give appropriate notions of correctness for basic systems. These include notions appropriate for systems that manage orphans, as well as notions for systems that do not manage orphans but do carry out concurrency control and recovery. These notions are taken from [11, 12, 4]. The spirit of our definitions is similar to that of the usual definition of "serializability" in the database literature. However, the usual notion does not take nesting or aborts into account.

We define correctness conditions for basic systems of a given type by relating their behaviors to those of a particular basic system of that type, the "serial system." The executions, schedules and behaviors of a serial system are called "serial executions," "serial schedules" and "serial behaviors," respectively. Serial systems are composed of transaction automata and a "serial database," which itself is the composition of a "serial scheduler" and objects. The transaction automata are identical to those in basic systems. The serial scheduler controls the order in which the transactions take steps and in which accesses to objects occur. It permits only one child of a transaction to run at a time. Thus, sibling transactions execute sequentially at every parent node in the transaction tree, so that transactions are run in a depth-first traversal of the tree. Also, the serial scheduler aborts a transaction only if it has not been created, and creates a transaction only if it has not been aborted. Thus, in a serial execution, sibling transactions execute sequentially and aborted transaction take no steps.

Objects in a serial system are quite simple. Since the serial scheduler guarantees that siblings execute sequentially, and that aborted transactions never take any steps, serial objects do not have to deal with concurrency or with failures. The serial objects serve as a specification of how objects should behave in the absence of concurrency and failures. (The serial objects serve the same purpose as the "serial specifications" in [25, 26].) A detailed description of serial systems may be found in the references [11, 12, 4, 13].

Now we give a definition that says that a sequence of actions "looks like" a serial behavior to a particular transaction. Namely, if  $\beta$  is a sequence of actions and  $T$  is a transaction name, we say that  $\beta$  is *serially correct for*  $T$  if there exists a serial behavior  $\gamma$  such that  $\gamma T = \beta T$ . In other words,  $T$  sees the same thing in  $\beta$  that it could see in some serial behavior.

Now we can define two notions of correctness for basic systems. First, we say that a basic system is *serially correct* if each of its finite<sup>13</sup> behaviors is serially correct for all transaction names.<sup>14</sup> The requirement that every transaction see a serial view is very strong. Without orphan management, in fact, systems may not meet this requirement. (This is true of all published concurrency control algorithms for nested transactions of which we are aware.) Instead, they provide a slightly weaker notion of correctness, namely that non-orphan transactions see serial views. More precisely, we say that a basic system is *serially correct for non-orphans* if each of its finite behaviors  $\beta$  is serially correct for all transaction names that are

---

<sup>13</sup>Serial correctness is stated in terms of finite behaviors because the corresponding property for infinite behaviors is not satisfied by locking algorithms, in the absence of extra assumptions [22].

<sup>14</sup>As discussed in [11], this definition of correctness allows different transactions in  $\beta$  to "see" different serial behaviors. However, correctness applies to the root transaction  $T_0$  as well, so the root must see the same results from the top-level transactions that it could see in some serial behavior.

not orphans in  $\beta$ . Orphans, however, can see arbitrary views.

The papers [11, 12, 4, 2] contain examples of basic systems that are serially correct for non-orphans. The basic system in [11, 12] uses exclusive locking for concurrency control and recovery,<sup>15</sup> while the systems in [4] use a more general commutativity-based locking strategy. The systems in [2] use timestamps for concurrency control and recovery.

The orphan management algorithms of this paper ensure that the systems that use them are serially correct. To ensure this, the orphan management algorithms rely on the basic database to ensure serial correctness for non-orphans; in fact, the algorithms work with any basic database that ensures serial correctness for non-orphans. In this sense, the orphan management algorithms and the concurrency control algorithms are independent. We prove a result of the following sort for each orphan management algorithm: if  $\beta$  is a behavior of the system with orphan management and  $T$  is a transaction name, then there exists a behavior  $\gamma$  of the underlying basic system such that  $\gamma T = \beta T$  and  $T$  is not an orphan in  $\gamma$ . In other words, the orphan management algorithms prevent transactions from “knowing” that they are orphans — everything a transaction sees is consistent with what it could see in the underlying basic system, in some execution in which it is not an orphan. These results imply that if the basic system is serially correct for non-orphans, then the corresponding system with orphan management is serially correct.

## 4. Information Flow

In this section, we define families of irreflexive partial orders, each of which models the information flow between events in behaviors of a basic system. These partial orders are used explicitly by the orphan management algorithms described later, in order to ensure that no transaction  $T$  ever obtains knowledge of the abort of any of its ancestors. If this is the case, then there will always be a “possible world” in which  $T$  is not an orphan and the interaction with  $T$  and its environment is unchanged.

### 4.1. Families of Affects Relations

If  $\beta$  is a sequence of basic actions,  $R$  is a binary relation on events in  $\beta$ , and  $\gamma$  is a subsequence of  $\beta$ , then we say that  $\gamma$  is *R-closed* in  $\beta$  if, whenever  $\gamma$  contains an event  $\pi$  in  $\beta$ , it also contains any event  $\phi$  such that  $(\phi, \pi) \in R$ .

Let  $B$  be a basic system, and let  $R = \{R_\beta\}$  be a family of relations, one for each sequence  $\beta$  of external actions of  $B$ . Then  $R$  is said to be a *family of affects relations* for  $B$  provided that the following conditions hold.

1. Each  $R_\beta$  is an irreflexive partial order on the events in  $\beta$  such that if  $(\phi, \pi) \in R_\beta$  then  $\phi$  precedes  $\pi$  in  $\beta$ .
2. If  $\gamma$  is a prefix of  $\beta$ , and  $\phi$  and  $\pi$  are in  $\gamma$ , then  $(\phi, \pi) \in R_\beta$  if and only if  $(\phi, \pi) \in R_\gamma$ .
3. If  $\beta$  is a behavior of  $B$  and  $\gamma$  is an  $R_\beta$ -closed subsequence of  $\beta$ , then  $\gamma$  is a behavior.

---

<sup>15</sup>There are some minor differences; for example, the completion and report actions are combined into single actions rather than treated as two separate actions.



of B.

4. Suppose that  $\beta$  is a behavior of B and  $(\phi, \pi) \in R_\beta$ , where  $\phi$  is an ABORT( $T'$ ) event, and  $\pi$  is a CREATE, a COMMIT, an ABORT, a REPORT\_COMMIT, a REPORT\_ABORT( $T$ ) for  $T \neq T'$ , an output of a non-access transaction, or a REQUEST\_COMMIT for an access. Then there is an event  $\psi$  between  $\phi$  and  $\pi$  in  $\beta$  such that  $(\phi, \psi) \in R_\beta$ , where  $\psi$  is related to  $\pi$  as follows:
  - a. If  $\pi = \text{CREATE}(T)$  then  $\psi = \text{REQUEST\_CREATE}(T)$ .
  - b. If  $\pi = \text{COMMIT}(T)$  then  $\psi = \text{REQUEST\_COMMIT}(T, v)$ .
  - c. If  $\pi = \text{REPORT\_COMMIT}(T, v)$  then  $\psi = \text{COMMIT}(T)$ .
  - d. If  $\pi = \text{REPORT\_ABORT}(T)$  then  $\psi = \text{ABORT}(T)$ .
  - e. If  $\pi = \text{ABORT}(T)$  then  $\psi = \text{REQUEST\_CREATE}(T)$ .
  - f. If  $\pi$  is an output of non-access transaction  $T$ , then  $\psi$  is an event of transaction  $T$ .
  - g. If  $\pi = \text{REQUEST\_COMMIT}(T, v)$  where  $T$  is an access to object  $X$ , then  $\text{object}(\psi) = X$ .<sup>16</sup>

The first two conditions are quite simple; the first says that each relation  $R_\beta$  describes a partial ordering consistent with the order in which events appear in  $\beta$ , and the second says that whether or not one event affects another is determined at the time the second event occurs. The third condition describes a sense in which the relation  $R_\beta$  captures all the dependency relationships between events. The condition implies that if  $\pi$  is not affected by  $\phi$  in some behavior  $\beta$ , then  $\pi$  cannot "know" that  $\phi$  occurred, since  $\pi$  could also have occurred in a different behavior in which  $\phi$  did not occur. The fourth condition is a technical condition that describes certain limitations on the pattern of information flow. It will be needed for our later proofs, and can be demonstrated for the examples in Section 9.

When the particular family  $R = \{R_\beta\}$  of affects relations is understood, we will often refer to each  $R_\beta$  as an "affects relation," and we will often say " $\phi$  affects  $\pi$  in  $\beta$ " to mean that  $(\phi, \pi) \in R_\beta$ .

#### 4.2. Families of Directly-Affects Relations

In many cases of interest, a family of affects relations can be conveniently described by a generating family of smaller relations. Thus, let B be a basic system and  $R' = \{R'_\beta\}$  be a family of relations, one for each sequence  $\beta$  of external actions of B. Then  $R'$  is said to be a *family of directly-affects relations* for B provided that there is a family  $R = \{R_\beta\}$  of affects relations for B such that for each  $\beta$ ,  $R_\beta$  is the transitive closure of  $R'_\beta$ . In this case, we say that  $R'$  *generates*  $R$ .

When the particular family  $R' = \{R'_\beta\}$  of directly-affects relations is understood, we will often refer to each  $R'_\beta$  as a "directly-affects relation"; also, we will often say that " $\phi$  directly affects  $\pi$  in  $\beta$ " to mean that  $(\phi, \pi) \in R'_\beta$ .

---

<sup>16</sup>Note that  $\psi$  can be either a serial or a non-serial event.

### 4.3. Using Families of Affects Relations to Describe Orphan Management Algorithms

The intuitive idea behind the orphan management algorithms is that they ensure that an event of a transaction  $T$  is never affected by the abort of an ancestor. Then we can show that every transaction gets a view it could get in a behavior in which it is not an orphan: we simply take the subsequence of the original behavior containing all events of  $T$  and all events that affect them in that behavior. The resulting sequence is a basic behavior, by the definition of a family of affects relations, and does not contain an abort for an ancestor of  $T$ , by construction.

The following example illustrates how an orphan can see an inconsistent state in a system without orphan management. Suppose that  $T$  is a transaction with children  $T_1$  and  $T_2$ , both of which are accesses to an object  $X$ . Consider the following scenario:  $T$  first accesses  $X$  through its child  $T_1$ .  $T$  then requests the creation of  $T_2$ , which will access  $X$  again. Furthermore, assume that  $T_2$  is requested only if  $T_1$  completes successfully, and that  $T_1$  modifies  $X$ . Now suppose that  $T$  aborts before  $T_2$  starts running at  $X$ , and  $X$  learns of the abort. If  $T_1$ 's modification of  $X$  is undone when  $X$  learns that  $T$  aborted, then  $T_2$  will not see the value for  $X$  that it expects, since  $T_2$  only runs if  $T_1$  has modified  $X$  successfully. This scenario is captured more precisely by the following fragment of a schedule of a generic system (generic systems are described in more detail in Section 9.1):

```
CREATE(T)
REQUEST_CREATE(T1)
CREATE(T1)
REQUEST_COMMIT(T1,v1)
COMMIT(T1)
REPORT_COMMIT(T1,v1)
REQUEST_CREATE(T2)
ABORT(T)
INFORM_ABORT_AT(X)OF(T)
CREATE(T2)
REQUEST_COMMIT(T2,v2)
...
```

(The INFORM\_ABORT event lets  $X$  know that  $T$  has aborted.) The family of affects relations for generic systems described in Section 9.1 ensures that the ABORT( $T$ ) event affects the INFORM\_ABORT\_AT( $X$ )OF( $T$ ) event, and that an event at an object is affected by all prior events at the object. Thus, by preventing  $T_2$  from running when its events would be affected by the abort of an ancestor, we can prevent it from knowing that it is an orphan.

## 5. Filtered Systems

The two orphan management algorithms analyzed in this paper use quite different techniques. However, each can be proved correct by showing that it implements the same abstract algorithm, described in this section.

For Sections 5 through 8 of this paper, we fix a particular but arbitrary basic system  $B$ , together with a family  $R$  of affects relations for  $B$  and a family  $R'$  of directly-affects relations for  $B$ , where  $R'$  generates  $R$ . We describe several algorithms that exploit the properties of the affects relations to manage orphans. In Section 9, we illustrate this general development by

describing two specific basic systems and their affects relations.

One way of ensuring that actions of a transaction  $T$  are never affected by the abort of an ancestor of  $T$  is to add preconditions to all the actions of the basic database to permit actions of  $T$  to occur only if they would not be affected in this way. It turns out, however, that this approach checks for orphans much more frequently than necessary. In this section we define another kind of system, called a "filtered system," that checks for orphans only when REQUEST\_COMMIT actions occur for access transactions. We then show that this is sufficient to ensure that transactions are never affected by the aborts of ancestors.

We construct a filtered system based on the given basic system  $B$  and the given family  $R$  of affects relations. The filtered system consists of the given transaction automata from  $B$  and a "filtered database automaton." The filtered database automaton is obtained by slightly modifying the basic database automaton; it "filters" REQUEST\_COMMIT actions of access transactions so that any transaction, orphan or not, sees a view it could see as a non-orphan in the basic system.

### 5.1. The Filtered Database

The filtered database is obtained via a simple transformation from the basic database. The only difference between the behaviors of the two databases is that the new database only allows a REQUEST\_COMMIT of an access to occur if it is not affected by the abort of an ancestor.

The filtered database has the same signature as the basic database. The state of the filtered database has two components, *basic\_state* and *history*, where *basic\_state* is a state of the basic database and *history* is a sequence of basic actions. Initial states of the filtered database are those with *basic\_state* equal to an initial state of the basic database and *history* equal to the empty sequence.

A triple  $(s', \pi, s)$  is a step of the filtered database if and only if the following conditions hold.

1.  $(s'.\text{basic\_state}, \pi, s.\text{basic\_state})$  is a step of the basic database.
2.  $s.\text{history} = s'.\text{history}\pi$  if  $\pi$  is a basic action.<sup>17</sup>
3.  $s.\text{history} = s'.\text{history}$  if  $\pi$  is not a basic action.
4. If  $\pi = \text{REQUEST\_COMMIT}(T, v)$  where  $T$  is an access to object  $X$ , and if  $T'$  is an ancestor of  $T$ , then no ABORT( $T'$ ) event affects an event  $\psi$  with  $\text{object}(\psi) = X$  in  $s'.\text{history}$ .

Thus, at the point where the REQUEST\_COMMIT of an access is about to occur, an explicit test is performed to verify that no preceding event at the same object is affected by the abort of any ancestor of the access.

**Lemma 6:** Let  $\beta$  be a finite schedule of the filtered database that can leave the filtered database in state  $s$ . Then  $s.\text{history} = \text{beh}(\beta)$ .

---

<sup>17</sup>Recall that internal actions of the basic database are not classified as basic actions.

## 5.2. The Filtered System

The *filtered system* is the composition of the transaction automata and the filtered database automaton. We call its executions, schedules and behaviors the *filtered executions*, *filtered schedules* and *filtered behaviors*, respectively.

**Lemma 7:** The filtered system implements the basic system.

**Proof:** The mapping  $f$  that assigns to each state  $s$  of the filtered system the singleton set  $f(s)$  consisting of  $s.basic\_state$  is easily seen to be a possibilities mapping. Proposition 3 implies the result.  $\square$

As described above, the filtered database performs an explicit test to ensure that the REQUEST\_COMMIT of an access is not affected by the abort of any ancestor. The following key lemma shows that this test actually guarantees more: that a similar property holds for all events.

**Lemma 8:** Let  $\beta$  be a filtered behavior and let  $T$  be any transaction name. Let  $\rho$  be an event in  $\beta$  such that  $transaction(\rho) = T$ . Then there is no ABORT( $T'$ ) event  $\phi$  such that  $(\phi, \rho) \in R_\beta$ , for any ancestor  $T'$  of  $T$ .

**Proof:** First note that Lemma 7 and Proposition 5 imply that  $serial(\beta)$  is basic database well-formed. The proof of the lemma is by induction on the length of  $\beta$ . If  $\beta$  is empty, the result clearly holds. Suppose  $\beta = \beta'\pi$ , and that the lemma holds for  $\beta'$ . From the restrictions on affects relations,  $R_\beta \subseteq R_{\beta'} \cup \{(\phi, \pi) \mid \phi \text{ is an action in } \beta'\}$ . Thus, by induction, it suffices to show that the lemma holds when  $\rho = \pi$ .

Suppose that the lemma does not hold, i.e., that  $\phi = \text{ABORT}(T')$  affects  $\pi$  in  $\beta$ , where  $transaction(\pi) = T$  and  $T'$  is an ancestor of  $T$ . We derive a contradiction. We consider cases.

1.  $T$  is a non-access and  $\pi$  is an output action of  $T$ .

Then by the fourth property of affects relations, there is an event  $\psi$  of  $T$  between  $\phi$  and  $\pi$  such that  $\phi$  affects  $\psi$  in  $\beta'$ . This contradicts the inductive hypothesis.

2.  $T$  is an access to object  $X$  and  $\pi$  is a REQUEST\_COMMIT for  $T$ .

Then by the fourth property of affects relations, there is an event  $\psi$  of object  $X$  between  $\phi$  and  $\pi$  in  $\beta'$  such that  $\phi$  affects  $\psi$  in  $\beta'$ . Then the precondition for  $\pi$  in the filtered database is violated, a contradiction.

3.  $\pi$  is CREATE( $T$ ).

Then by the fourth property of affects relations, there is a REQUEST\_CREATE( $T$ ) event  $\psi$  between  $\phi$  and  $\pi$  such that  $\phi$  affects  $\psi$  in  $\beta'$ . Since REQUEST\_CREATE( $T$ ) is an action of  $parent(T)$ , the inductive hypothesis implies that  $T'$  is not an ancestor of  $parent(T)$ . The only possibility is that  $T' = T$ , which implies that ABORT( $T$ ) precedes REQUEST\_CREATE( $T$ ) in  $\beta$ . But this implies that  $serial(\beta)$  is not basic database well-formed, a contradiction.

4.  $T$  is a non-access and  $\pi$  is REPORT\_COMMIT( $T'', v$ ), where  $T''$  is a child of  $T$ .

Then  $T'$  is an ancestor of  $T''$ . By the fourth property of affects relations, there is a REQUEST\_COMMIT( $T'', v$ ) event  $\psi$  between  $\phi$  and  $\pi$  such that  $\phi$  affects  $\psi$  in  $\beta'$ . Since  $transaction(\psi) = T''$ , this contradicts the inductive hypothesis.

5.  $T$  is a non-access and  $\pi$  is  $\text{REPORT\_ABORT}(T'')$ , where  $T''$  is a child of  $T$ .

By the fourth property of affects relations, there is a  $\text{REQUEST\_CREATE}(T'', v)$  event  $\psi$  between  $\phi$  and  $\pi$  such that  $\phi$  affects  $\psi$  in  $\beta'$ . Since  $\text{transaction}(\psi) = T$ , this contradicts the inductive hypothesis.

6.  $\pi$  is a non-serial basic action.

Then  $\text{transaction}(\pi)$  is undefined, a contradiction.

□

### 5.3. Simulation of the Basic System by the Filtered System

The following theorem is the key result of this paper. It shows that the filtered system ensures that every transaction gets a view it could get in the basic system when it is not an orphan. (Formally, a transaction  $T$ 's "view" in a behavior  $\beta$  is its local behavior,  $\beta|T$ .) In other words, an orphan cannot discover that it is an orphan, since the view it sees is consistent with its not being an orphan. This is the basic correctness property for the orphan management algorithms.

**Theorem 9:** Let  $\beta$  be a filtered behavior and let  $T$  be a transaction name. Then there exists a basic behavior  $\gamma$  such that  $T$  is not an orphan in  $\gamma$  and  $\gamma|T = \beta|T$ .

**Proof:** Let  $\gamma$  be the subsequence of  $\beta$  containing all actions  $\pi$  such that  $\text{transaction}(\pi) = T$ , and all other actions  $\phi$  that affect, in  $\beta$ , some action whose transaction is  $T$ . Since  $R_\beta$  is a transitive relation,  $\gamma$  is  $R_\beta$ -closed in  $\beta$ . By the definition of a family of affects relations,  $\gamma$  is a basic behavior. It suffices to show that there is no ancestor  $T'$  of  $T$  for which  $\text{ABORT}(T')$  occurs in  $\gamma$ . Suppose not; i.e., there exists an ancestor  $T'$  of  $T$  for which  $\text{ABORT}(T')$  occurs in  $\gamma$ . Then by the construction of  $\gamma$ ,  $\beta$  contains an event  $\pi$  of  $T$  such that an  $\text{ABORT}(T')$  event  $\phi$  affects  $\pi$  in  $\beta$ . By Lemma 8, this is impossible. □

We obtain an important corollary of Theorem 9.

**Corollary 10:** If the basic system is serially correct for non-orphans, then the filtered system is serially correct.<sup>18</sup>

**Proof:** Let  $\beta$  be a filtered behavior and let  $T$  be any transaction name. Theorem 9 yields a behavior  $\delta$  of the basic system such that  $T$  is not an orphan in  $\delta$  and  $\delta|T = \beta|T$ . Since the basic system is serially correct for non-orphans, there is a serial behavior  $\gamma$  with  $\gamma|T = \delta|T$ ; this is equal to  $\beta|T$ , as needed. □

At first it might seem somewhat surprising that it is enough to prevent the  $\text{REQUEST\_COMMIT}$  events of orphan accesses to ensure serial correctness for all orphans. The reason it is not necessary to filter other actions is because of the assumptions about affects relations. Essentially, these assumptions indicate that an event, say  $\text{CREATE}(T)$ , cannot "know" about an  $\text{ABORT}(T')$  event unless some earlier event, in this case  $\text{REQUEST\_CREATE}(T)$ , already "knows" about the  $\text{ABORT}(T')$ . (The assumption about affects relations that an affects-closed subsequence of a behavior is itself a behavior implies that if  $\pi$  is not affected by  $\phi$ ,  $\pi$  cannot know that  $\phi$  has occurred, since there is a behavior of the system in which  $\pi$  occurs and  $\phi$  does not.) If we assume inductively that

<sup>18</sup>It should be easy to see that the filtered system is also a basic system, and so the notion of serial correctness has been defined for filtered systems. Similar comments hold for the rest of the systems described in this paper.

REQUEST\_CREATE(T) does not know about the abort of an ancestor of T, then the properties assumed for affects relations guarantee that CREATE(T) will not know about it either.

The assumptions about affects relations derive from the restricted communication patterns in typical systems: a non-access transaction receives information from its parent when it is created, and from its children when they report, but not from any other source. Access transactions may receive information from other accesses (e.g., accesses to the same object share state), but can only affect non-accesses through a REPORT\_COMMIT event, which must be preceded by a REQUEST\_COMMIT. As long as T does not receive reports from any accesses that "know" that its ancestor has aborted, T cannot observe a state that depends on the abort. In effect, by preventing REQUEST\_COMMIT actions for orphan accesses, we isolate orphan transactions from the objects, ensuring that an orphan transaction never sees that it is an orphan.

Sometimes we want to be explicit about the dependency of the filtered system on the particular basic system and family of affects relations from which it is derived. Thus, we restate the corollary above in a form that exhibits this dependency. If B is a basic system and R is a family of affects relations for B then let Filtered(B,R) be the corresponding filtered system; it is composed of the same transaction automata and the filtered database that corresponds to the given basic database.

**Corollary 11:** Let B be a basic system and R a family of affects relations for B. If B is serially correct for non-orphans, then Filtered(B,R) is serially correct.

## 6. Argus Systems

In this section we analyze the orphan management algorithm used in the Argus system [8, 9]. We describe the algorithm by defining an *Argus database* that describes in formal terms the algorithm discussed in [9]. As with the filtered database, the Argus database is obtained from the basic database via a simple construction. We then define the Argus system, which is composed of transactions and an Argus database, and show that the Argus system implements the filtered system. Thus, if the filtered system is serially correct, so is the corresponding Argus system.

### 6.1. The Argus Database

The filtered database uses global knowledge of the entire history of actions to filter the REQUEST\_COMMIT actions of access transactions. This kind of global knowledge is not practical in a distributed system. Thus, the Argus algorithm makes use of local knowledge about the aborts that have occurred. To ensure that the REQUEST\_COMMIT of an access is not affected by the abort of an ancestor, the Argus algorithm keeps track of the aborts "known" by each event that occurs, and propagates this knowledge from an event to any later events that it affects.

In the actual Argus system, knowledge about aborts is propagated in messages sent over the network; we model this formally by propagating knowledge from an event to every event that it *directly affects*. A poor choice of a directly-affects relation could make the algorithm we describe here hard to implement. However, if one event  $\phi$  directly affects another event  $\pi$  only if the two events occur at the same site, or if a message is sent from  $\phi$ 's site to  $\pi$ 's after  $\phi$  occurs and before  $\pi$  occurs, then it is straightforward to implement the algorithm using information

available locally at each site, by transmitting the information about aborted transactions on messages sent between sites.

The construction of the Argus database makes use of the given family  $R'$  of directly-affects relations for  $B$ . The Argus database has the same signature as the basic database. The state of the Argus database has three components: *basic\_state*, *history* and *known\_aborts*, where *basic\_state* is a state of the basic database, *history* is a sequence of basic actions, and *known\_aborts* is a partial mapping from basic events to sets of transactions. This mapping records the transactions whose aborts affect each event that has occurred. The set  $\text{known\_aborts}(\pi)$  may actually include more transactions than those whose aborts affect  $\pi$ . By adding more aborted transactions to this set, an implementation would restrict the behavior of orphans further than is strictly necessary to ensure the correctness conditions.<sup>19</sup> In Argus, for example, each event occurs at some physical node of the network, and each node manages a single set of aborted transactions for the entire set of events that occur at that node. In this case,  $\text{known\_aborts}(\pi)$  includes at least the transactions whose aborts affect  $\pi$ , as well as those transactions whose aborts affect any other event that has occurred at the same physical node as  $\pi$ .

Initial states of the Argus database are those with *basic\_state* equal to an initial state of the basic database, *history* equal to the empty sequence, and *known\_aborts* everywhere undefined. A triple  $(s', \pi, s)$  is a step of the Argus database if and only if the following conditions hold.

1.  $(s'.\text{basic\_state}, \pi, s.\text{basic\_state})$  is a step of the basic database.
2.  $s.\text{history} = s'.\text{history}\pi$  if  $\pi$  is a basic action.
3.  $s.\text{history} = s'.\text{history}$  if  $\pi$  is not a basic action.
4. If  $\pi$  is a basic action and  $(\phi, \pi) \in R'_{s.\text{history}}$  then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$ .
5. If  $\phi \neq \pi$ , then  $s.\text{known\_aborts}(\phi) = s'.\text{known\_aborts}(\phi)$ .
6. If  $\pi$  is a basic action and  $\phi$  is an  $\text{ABORT}(T)$  event in  $s'.\text{history}$  such that  $(\phi, \pi) \in R'_{s'.\text{history}}$  then  $T \in s.\text{known\_aborts}(\pi)$ .
7. If  $\pi$  is a  $\text{REQUEST\_COMMIT}(T, v)$  action for an access  $T$  to object  $X$ , then there is no ancestor of  $T$  in  $s'.\text{known\_aborts}(\phi)$ , for any event  $\phi$  of object  $X$  in  $s'.\text{history}$ .

There are two significant differences between the Argus database and the basic database. First, the effects for each action  $\pi$  in the Argus database require  $s.\text{known\_aborts}(\pi)$  to include  $s'.\text{known\_aborts}(\phi)$  for each  $\phi$  that directly affects  $\pi$ . In addition, an event  $\pi$  that is directly affected by  $\text{ABORT}(T)$  requires  $T$  to be in  $\text{known\_aborts}(\pi)$ . As Lemma 16 below shows, these constraints are enough to ensure that  $s.\text{known\_aborts}(\pi)$  contains  $T$  whenever  $\text{ABORT}(T)$  affects  $\pi$ .

Second, the precondition for the  $\text{REQUEST\_COMMIT}$  of an access to  $X$  permits the event to

---

<sup>19</sup>In fact, our description does not require the *known\_aborts* set for an event to contain only aborted transactions. Adding non-aborted transactions to the set might cause non-orphans to block, but will otherwise not result in incorrect behavior. One might like to prove that a system does not treat a non-orphan as an orphan; as with other liveness properties, we do not address this issue here.

occur only if the access does not "know about" the abort of an ancestor, i.e., no ancestor is in  $s'.\text{known\_aborts}(\phi)$  for any event  $\phi$  of object  $X$  in  $s'.\text{history}$ . As Lemma 17 below shows, this is enough to ensure that every Argus behavior is a filtered behavior.

The `known_aborts` mapping models the distributed information maintained by the Argus algorithm to keep track of actions that abort. However, rather than modeling nodes directly and keeping the information on a per-node basis as is done in the actual algorithm, we maintain the information for each event, propagating it whenever one event directly affects another.

The `known_aborts` component is managed so as to insure that at least the minimum amount of necessary information is propagated at each step. An implementation is permitted to propagate more than the minimum; for instance, an implementation might keep track of the `known_aborts` mapping at a coarser granularity. (By maintaining the `known_aborts` mapping on a per-node basis, the implementation of the Argus algorithm in the current Argus prototype follows this strategy.) In describing the algorithm, we have tried to focus on the behavior necessary for correctness, and to avoid constraining an implementation any more than necessary.

Notice also that the Argus database does not put any upper limit on what goes into the `known_aborts` mapping. For example, as described above it is permissible for `known_aborts( $\pi$ )` to contain a transaction that has not aborted. It would be easy (and intuitively appealing) to add a requirement that `known_aborts( $\pi$ )` only includes aborted transactions, but this is not necessary to prove that the algorithm prevents all orphans from seeing inconsistent views. To prove other properties, such as that the algorithm only detects real orphans, we would need to add additional requirements such as the one just mentioned. We will not attempt to state or prove such properties in this paper; the property just described is a special case of more general liveness properties, which are an appropriate subject of further research.

Finally, while the Argus database is distinguished from the filtered database by maintaining information about ABORT actions on a more local basis, we have kept the state component `s.history`, which maintains global knowledge of the past behavior of the system. A practical implementation of this algorithm would maintain less voluminous history information in a distributed fashion. Examination of the additional preconditions imposed by the Argus database on any action  $\pi$  reveals that `s.history` is used to determine the events  $\phi$  that directly affect  $\pi$ , and when  $\pi$  is a `REQUEST_COMMIT(T,v)` action for an access  $T$  to object  $X$ , to determine the events that precede  $\pi$  at  $X$ . The details of efficient maintenance of sufficient history information are dependent on the particular basic database and distribution scheme, and are not addressed further in this paper.

## 6.2. The Argus System

The *Argus system* is the composition of transactions and the Argus database. Executions, schedules and behaviors of the Argus system are called *Argus executions*, *Argus schedules* and *Argus behaviors*, respectively.

**Lemma 12:** The Argus system implements the basic system.

**Proof:** The proof is similar to that of Lemma 7. □

**Lemma 13:** Let  $\beta$  be a finite Argus behavior that can leave the Argus database in state  $s$ . Then `s.known_aborts( $\pi$ )` is defined if and only if  $\pi$  is an event in  $\beta$ .



The following lemma says that each `known_aborts` set is defined at most once during an Argus execution.

**Lemma 14:** Let  $\beta'\beta$  be a finite Argus schedule, where  $\beta'$  can leave the Argus database in state  $s'$  and  $(s',\beta,s)$  is an extended step of the Argus database. If  $s'.\text{known\_aborts}(\pi)$  is defined, then  $s'.\text{known\_aborts}(\pi) = s.\text{known\_aborts}(\pi)$ .

The next lemma says that the `known_aborts` set for an event  $\pi$  includes all events that directly affect  $\pi$ , and that the `known_aborts` set for  $\pi$  includes  $T$  if  $\pi$  is directly affected by an `ABORT(T)` event.

**Lemma 15:** Let  $\beta$  be a finite Argus behavior that can leave the Argus database in state  $s$ .

1. If  $\phi$  and  $\pi$  are events in  $\beta$  such that  $\phi$  directly affects  $\pi$  in  $\beta$ , then  $s.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$ .
2. If  $\pi$  is directly affected by an `ABORT(T)` event in  $\beta$  then  $T \in s.\text{known\_aborts}(\pi)$ .

**Proof:** Immediate by the definition of the Argus database steps and Lemmas 13 and 14.  $\square$

The next lemma is the key to the proof of correctness for the Argus system: it says that the `known_aborts` set for an event  $\pi$  includes all transactions  $T$  such that an `ABORT(T)` event affects  $\pi$ . In other words, the Argus database propagates enough information about aborts so that every event  $\pi$  "knows about" (stores in  $s.\text{known\_aborts}(\pi)$ ) every abort that affects it.

**Lemma 16:** Let  $\beta$  be a finite Argus behavior that can leave the Argus database in state  $s$ . If  $\phi$  and  $\pi$  are events in  $\beta$  such that  $\phi$  affects  $\pi$  in  $\beta$  and  $\phi$  is an `ABORT(T)` event, then  $T \in s.\text{known\_aborts}(\pi)$ .

**Proof:** The proof proceeds by induction on the length of the chain in the directly-affects relation by which  $\phi$  affects  $\pi$  in  $\beta$ . If the length of the chain is 1, then  $\phi$  directly affects  $\pi$  in  $\beta$ . Then Lemma 15 implies that  $T \in s.\text{known\_aborts}(\pi)$ .

Now suppose that the length of the chain is  $k+1$ , where  $k \geq 1$ . Then there is an event  $\psi$  in  $\beta$  such that  $\phi$  affects  $\psi$  in  $\beta$  by a chain of length at most  $k$  and  $\psi$  directly affects  $\pi$  in  $\beta$ . By inductive hypothesis,  $T \in s.\text{known\_aborts}(\psi)$ . Lemma 15 implies that  $s.\text{known\_aborts}(\psi) \subseteq s.\text{known\_aborts}(\pi)$ , so that  $T \in s.\text{known\_aborts}(\pi)$ .  $\square$

### 6.3. Simulation of the Basic System by the Argus System

The following lemma shows that the information in `known_aborts`, combined with the precondition on `REQUEST_COMMIT` actions for accesses, is enough to ensure that the Argus system implements the filtered system.

**Lemma 17:** The Argus system implements the filtered system.

**Proof:** We define a mapping  $f$  that assigns to each state  $s$  of the Argus system the singleton set  $f(s)$  consisting of the state of the filtered system that is the same except for the omission of the `s.known_aborts` component of the Argus database state. We must show that  $f$  is a possibilities mapping. Condition 1. is easy to check. For condition 2., suppose that  $s'$  is a reachable state of the Argus system,  $t' \in f(s')$  is a reachable state of the filtered system, and  $(s',\pi,s)$  is a step of the Argus system. The

only interesting case to check is when  $\pi = \text{REQUEST\_COMMIT}(T, v)$ , where  $T$  is an access to an object  $X$ . In this case, we claim that  $(t', \pi, t)$  is a step of the filtered system, where  $t$  is the single element of  $f(s)$ .

To show that  $(t', \pi, t)$  is a step of the filtered system, we must show the four conditions defining these steps. The first three are immediate from the definition of the steps of the Argus database. To see the fourth, suppose that  $T'$  is an ancestor of  $T$ . We must show that no  $\text{ABORT}(T')$  event affects an event of object  $X$  in  $t'$ .history. Suppose the contrary, that an  $\text{ABORT}(T')$  event  $\phi$  affects an event  $\psi$  of object  $X$  in  $t'$ .history. Then  $\phi$  also affects  $\psi$  in  $s'$ .history. By Lemma 16,  $T' \in s.\text{known\_aborts}(\psi)$ . But this violates the precondition for  $\pi$  in the Argus database. Therefore,  $\pi$  is enabled in  $t'$ .  $\square$

The following theorem shows that Argus systems, like filtered systems, ensure that every non-access transaction gets a view it could get in an execution in which it is not an orphan.

**Theorem 18:** . Let  $\beta$  be an Argus behavior and let  $T$  be a transaction name. Then there exists a basic behavior  $\gamma$  such that  $T$  is not an orphan in  $\gamma$  and  $\gamma|T = \beta|T$ .

**Proof:** Immediate by Lemma 17 and Theorem 9.  $\square$

As for the filtered system, we obtain an important corollary about serial correctness.

**Corollary 19:** If the basic system is serially correct for non-orphans, then the Argus system is serially correct.

Again, we give a version of the preceding corollary in which the dependency of the Argus system on the basic system is made explicit. If  $B$  is a basic system and  $R'$  is a family of directly-affects relations for  $B$ , then let  $\text{Argus}(B, R')$  be the corresponding Argus system; it is composed of the same transaction automata and the Argus database that is constructed from the given basic database, using the given family of relations  $R'$ .

**Corollary 20:** Suppose  $B$  is a basic system and  $R'$  is a family of directly-affects relations for  $B$ . If  $B$  is serially correct for non-orphans, then  $\text{Argus}(B, R')$  is serially correct.

## 7. Strictly Filtered Systems

The orphan management algorithm described in [16] actually ensures a stronger property than does the Argus algorithm. It ensures that  $\text{REQUEST\_COMMIT}$  can never occur for an orphan access, whereas the Argus algorithm merely ensures that no such  $\text{REQUEST\_COMMIT}$  can occur if the access can "observe" that it is an orphan. In this section we define the "strictly filtered database," which allows a  $\text{REQUEST\_COMMIT}$  to occur for an access only if no ancestor has aborted. (Compare this to the filtered database, which allows an access to  $\text{REQUEST\_COMMIT}$  if an ancestor has aborted as long as the access is not affected by the abort.) We then define the strictly filtered system, which is composed of transactions and the strictly filtered controller, and show that the strictly filtered system implements the filtered system. In the next section we will describe formally the algorithm from [16] and show that it implements the strictly filtered system.

### 7.1. The Strictly Filtered Database

The strictly filtered database is similar to the filtered database; it has the same actions, and the same states. A triple  $(s', \pi, s)$  is a step of the strictly filtered database if and only if the following conditions hold.

1.  $(s'.\text{basic\_state}, \pi, s.\text{basic\_state})$  is a step of the basic database.
2.  $s.\text{history} = s'.\text{history}\pi$  if  $\pi$  is a basic action.
3.  $s.\text{history} = s'.\text{history}$  if  $\pi$  is not a basic action.
4. If  $\pi = \text{REQUEST\_COMMIT}(T, v)$  where  $T$  is an access, and  $T'$  is an ancestor of  $T$ , then no  $\text{ABORT}(T')$  event occurs in  $s'.\text{history}$ .

Thus, at the point where the  $\text{REQUEST\_COMMIT}$  of an access is about to occur, an explicit test is performed to verify that there is no preceding abort of any ancestor of the access.

### 7.2. The Strictly Filtered System

The *strictly filtered system* is the composition of transactions and the strictly filtered database. Executions, schedules and behaviors of the strictly filtered system are *strictly filtered executions*, *schedules* and *behaviors*, respectively.

**Lemma 21:** The strictly filtered system implements the basic system.

**Proof:** The proof is similar to that of Lemma 7. □

### 7.3. Simulation of the Basic System by the Strictly Filtered System

**Lemma 22:** The strictly filtered system implements the filtered system.

**Proof:** We define a mapping  $f$  that assigns to each state  $s$  of the strictly filtered system the singleton set  $f(s)$  that consists of the same state. We must show that  $f$  is a possibilities mapping. Condition 1. is easy to check. For condition 2., suppose that  $s'$  is a reachable state of the strictly filtered system,  $t' \in f(s')$  is a reachable state of the filtered system, and  $(s', \pi, s)$  is a step of the strictly filtered system. As before, the only interesting case to check is when  $\pi = \text{REQUEST\_COMMIT}(T, v)$ , where  $T$  is an access to an object  $X$ . In this case, we claim that  $(t', \pi, t)$  is a step of the filtered system, where  $t$  is the unique element of  $f(s)$ .

To show that  $(t', \pi, t)$  is a step of the filtered system, we must show the four conditions defining these steps. The first three are immediate from the definition of the steps of the strictly filtered database. To see the fourth, suppose that  $T'$  is an ancestor of  $T$ . We must show that no  $\text{ABORT}(T')$  event affects an event of object  $X$  in  $t'.\text{history}$ . But  $t'.\text{history} = s'.\text{history}$ , and by the preconditions for  $\pi$  in the strictly filtered database, no  $\text{ABORT}(T')$  event occurs in  $s'.\text{history}$ . Therefore, no  $\text{ABORT}(T')$  event affects an event of object  $X$  in  $t'.\text{history}$ . Thus,  $\pi$  is enabled in  $t'$ . □

Strictly filtered systems, like filtered systems and Argus systems, prevent orphans from discovering that they are orphans.

**Theorem 23:** Let  $\beta$  be a strictly filtered behavior and let  $T$  be a transaction name. Then there exists a basic behavior  $\gamma$  such that  $T$  is not an orphan in  $\gamma$  and  $\gamma \upharpoonright T = \beta \upharpoonright T$ .

**Proof:** Immediate by Lemma 22 and Theorem 9. □

**Corollary 24:** If the basic system is serially correct for non-orphans, then the strictly filtered system is serially correct.

If  $B$  is a basic system, let  $\text{Strictly-Filtered}(B)$  be the corresponding strictly filtered system; it is composed of the same transaction automata and the strictly filtered database that corresponds to the given basic database.

**Corollary 25:** Suppose that  $B$  is a basic system. If  $B$  is serially correct for non-orphans then  $\text{Strictly-Filtered}(B)$  is serially correct.

## 8. Clocked Systems

In this section we describe formally the orphan management algorithm from [16]. (The algorithm described here actually generalizes the one described in [16]. The algorithm in [16] delays commits of transactions, as well as aborts, while the algorithm described here delays only aborts. A similar generalization is described in [17].) We do this by defining the "clocked database," which uses a global clock to ensure that transactions do not abort until all their descendant accesses have stopped running. We then define the clocked system, which is composed of transactions and the clocked database. Finally, we show that the clocked system implements the strictly filtered system, and thus simulates the basic system in the same way as the previously mentioned systems do.

### 8.1. The Clocked Database

The clocked database maintains a quiesce time for each access transaction and a release time for every transaction. An access transaction is allowed to `REQUEST_COMMIT` only if its quiesce time has not passed. Release times are chosen so that once a transaction's release time is reached, all its descendant accesses have quiesced. A transaction is allowed to abort only if its release time has passed. This ensures that, after a transaction aborts, none of its descendant accesses will request to commit.

If quiesce and release times are fixed in advance, some transactions may be forced to abort unnecessarily as their quiesce times expire, and aborts may need to be delayed until release times are reached. It is possible to obtain extra flexibility by providing actions in the clocked database for adjusting quiesce and release times.

The action signature of the clocked database is the same as that of the basic database, except that the clocked database has three additional kinds of internal actions. The new actions are:

Internal Actions:

`TICK`

`ADJUST_QUIESCE(T)`,  $T$  an access

`ADJUST_RELEASE(T)`,  $T$  any transaction

The `TICK` action advances the clock, while the two `ADJUST` actions adjust quiesce and release times. By adjusting the quiesce time for a transaction to be later than its current value, we can extend the time during which a transaction is allowed to run. Similarly, by adjusting the release time for a transaction to be earlier than its current value, we can allow a transaction to abort without waiting as long as would otherwise be necessary.

The state of the clocked database consists of components *basic\_state*, *aborted*, *clock*, *quiesce*, and *release*. Here, *basic\_state* is a state of the basic database, initialized at an initial state of the basic database. The component *aborted* is a set of transactions, initially empty. The component *clock* is a real number, initialized arbitrarily. The component *quiesce* is a total mapping from access transaction names to real numbers, and the component *release* is a total mapping from all transaction names to real numbers. The initial values of *quiesce* and *release* are arbitrary, subject to the following condition: for all transaction names  $T$  and  $T'$ , where  $T$  is an access and  $T'$  is an ancestor of  $T$ ,  $quiesce(T) \leq release(T')$ .

A triple  $(s', \pi, s)$  is a step of the clocked database if and only if the following conditions hold.

1. If  $\pi$  is an action of the basic system, then  $(s'.basic\_state, \pi, s.basic\_state)$  is a step of the basic database.
2. If  $\pi$  is a TICK, ADJUST\_QUIESCE or ADJUST\_RELEASE action, then  $s.basic\_state = s'.basic\_state$ .
3. If  $\pi = ABORT(T)$  then
  - a.  $s.aborted = s'.aborted \cup \{T\}$ .
  - b.  $s'.release(T) \leq s'.clock$ .
4. If  $\pi$  is not an abort action, then  $s.aborted = s'.aborted$ .
5. If  $\pi = REQUEST\_COMMIT(T, v)$  where  $T$  is an access, then  $s'.clock < s'.quiesce(T)$ .
6. If  $\pi = TICK$  then  $s'.clock < s.clock$ .
7. If  $\pi$  is not a TICK action, then  $s.clock = s'.clock$ .
8. If  $\pi = ADJUST\_RELEASE(T)$  then
  - a. if  $T \in s'.aborted$  then  $s.release(T) \leq s'.clock$ ,
  - b.  $s'.quiesce(T') \leq s.release(T)$  for all  $T' \in descendants(T) \cap accesses$ , and
  - c.  $s.release(T') = s'.release(T')$  for all  $T' \neq T$ .
9. If  $\pi$  is not an ADJUST\_RELEASE action, then  $s.release = s'.release$ .
10. If  $\pi = ADJUST\_QUIESCE(T)$  then
  - a.  $s.quiesce(T) \leq s'.release(T')$  for all  $T' \in ancestors(T)$ , and
  - b.  $s.quiesce(T') = s'.quiesce(T')$  for all  $T' \neq T$ .
11. If  $\pi$  is not an ADJUST\_QUIESCE action, then  $s.quiesce = s'.quiesce$ .

**Lemma 26:** Let  $\beta$  be a finite schedule of the clocked database that can leave the clocked database in state  $s$ .

1. If  $T \in s.aborted$  then  $s.release(T) \leq s.clock$ .
2. For all accesses  $T$  and all ancestors  $T'$  of  $T$ ,  $s.quiesce(T) \leq s.release(T')$ .

**Proof:** Straightforward by induction. □

## 8.2. The Clocked System

The *clocked system* is the composition of transactions and the clocked database. External actions of the clocked system are called *clocked actions*. Executions, schedules and behaviors of a clocked system are called *clocked executions*, *schedules* and *behaviors*, respectively.

**Lemma 27:** The clocked system implements the basic system.

**Proof:** The proof is similar to that of Lemma 7. □

## 8.3. Simulation of the Basic System by the Clocked System

**Lemma 28:** The clocked system implements the strictly filtered system.

**Proof:** We define a mapping  $f$  that assigns to each state  $s$  of the clocked system the set  $f(s)$  of states  $t$  of the strictly filtered system such that  $t.\text{basic\_state} = s.\text{basic\_state}$  and  $t.\text{history}$  is a sequence of basic actions in which the set of transaction names  $T$  for which  $\text{ABORT}(T)$  occurs in  $t.\text{history}$  is exactly  $s.\text{aborted}$ . We must show that  $f$  is a possibilities mapping. Condition 1. is easy to check. For condition 2., suppose that  $s'$  is a reachable state of the clocked system,  $t' \in f(s')$  is a reachable state of the strictly filtered system and  $(s', \pi, s)$  is a step of the clocked system.

There are two interesting cases to check: where  $\pi = \text{ABORT}(T)$  and where  $\pi = \text{REQUEST\_COMMIT}(T, v)$  for an access  $T$ . In either case, we claim that  $(t', \pi, t)$  is a step of the strictly filtered system, where  $t$  is the state of the strictly filtered system in which  $t.\text{basic\_state} = s.\text{basic\_state}$  and  $t.\text{history} = t'.\text{history}\pi$ , and we also claim that  $t \in f(s)$ .

If  $\pi = \text{ABORT}(T)$ , it is easy to see that  $(t', \pi, t)$  is a step of the strictly filtered system. To show  $t \in f(s)$ , note that since  $t' \in f(s')$ , the set of transaction names  $U$  for which  $\text{ABORT}(U)$  occurs in  $t'.\text{history}$  is exactly  $s'.\text{aborted}$ . Then the set of transaction names with aborts in  $t.\text{history}$  is exactly  $s'.\text{aborted} \cup \{T\}$ , which is equal to  $s.\text{aborted}$ . Thus,  $t \in f(s)$ .

If  $\pi = \text{REQUEST\_COMMIT}(T, v)$ , where  $T$  is an access, then it is easy to see the first three conditions of the definition of strictly filtered database steps. For the fourth condition, we must show that if  $T'$  is an ancestor of  $T$ , then no  $\text{ABORT}(T')$  event occurs in  $t'.\text{history}$ . So suppose the contrary, that  $T'$  is an ancestor of  $T$  and  $\text{ABORT}(T')$  occurs in  $t'.\text{history}$ . Since  $t' \in f(s')$ , we have  $T' \in s'.\text{aborted}$ . Since  $\pi$  is enabled in  $s'$ ,  $s'.\text{clock} < s'.\text{quiesce}(T)$ . Lemma 26 implies that  $s'.\text{release}(T') \leq s'.\text{clock}$  and also that  $s'.\text{quiesce}(T) \leq s'.\text{release}(T')$ . Thus,  $s'.\text{quiesce}(T) \leq s'.\text{clock}$ , a contradiction. It follows that  $(t', \pi, t)$  is a step of the strictly filtered system.

Since  $t' \in f(s')$ , the set of transaction names  $U$  for which  $\text{ABORT}(U)$  occurs in  $t'.\text{history}$  is exactly  $s'.\text{aborted}$ . Then the set of transaction names with aborts in  $t.\text{history}$  is exactly  $s'.\text{aborted} = s.\text{aborted}$ . Thus,  $t \in f(s)$ . □

**Theorem 29:** Let  $\beta$  be a clocked behavior and let  $T$  be a transaction name. Then there exists a basic behavior  $\gamma$  such that  $T$  is not an orphan in  $\gamma$  and  $\gamma|T = \beta|T$ .

**Proof:** By Lemma 28 and Theorem 23. □

**Corollary 30:** If the basic system is serially correct for non-orphans, then the clocked system is serially correct.

If  $B$  is a basic system, then let  $\text{Clocked}(B)$  be the corresponding clocked system; it is

composed of the same transaction automata and the clocked database of the appropriate type.

**Corollary 31:** Suppose  $B$  is a basic system. If  $B$  is serially correct for non-orphans, then  $\text{Clocked}(B)$  is serially correct.

The algorithm described here uses a single physical clock to detect and eliminate orphans. The algorithm can be adapted to work with distributed, loosely synchronized physical clocks, or with logical clocks (e.g., see [17]). The adapted algorithms can be described and analyzed in a manner similar to that used for the Argus algorithm.

## 9. Examples

In this section, we describe two important kinds of basic systems, together with a family of affects relations (and a generating family of directly-affects relations) for each of them. The first kind of system is a "generic system." It is suitable for modeling locking algorithms and has been studied in [4]. The second is a "pseudotime system." It is suitable for modeling timestamp algorithms and has been studied in [2].

### 9.1. Generic Systems

A generic system consists of a collection of transaction automata, one for each non-access transaction name, a collection of "generic object automata," one for each object name, and a single "generic controller automaton." The interactions between the components are as follows.

The transaction interface is exactly as before. The generic object automaton for  $X$  has  $\text{CREATE}(T)$  input actions and  $\text{REQUEST\_COMMIT}(T,v)$  output actions for each access  $T$  to  $X$  and each return value  $v$ . It also has  $\text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$  and  $\text{INFORM\_ABORT\_AT}(X)\text{OF}(T)$  input actions for each transaction  $T$ ; these actions inform the object  $X$  of the fates (commit or abort) of completed transactions. The object uses this information in carrying out concurrency control and recovery; for example, an  $\text{INFORM\_ABORT}$  of  $T$  might cause the object to release locks held by  $T$ .

The external actions of the generic controller are similar to those required of all basic databases: it has the inputs and outputs required of a basic database, except that  $\text{REQUEST\_COMMIT}$  actions for access transactions are inputs to the generic controller; in addition, it has  $\text{INFORM\_COMMIT}$  and  $\text{INFORM\_ABORT}$  actions as outputs.

We model the generic controller as a specific automaton, particular to the system type. The object automata, however, like the transaction automata, are only partially specified. Their signature is as described above, and they are constrained to preserve an appropriately defined well-formedness property. Otherwise, they are unconstrained. In particular, the semantics of their operations is immaterial to our discussion. In this paper, we are concerned only with whether the entire generic database is serially correct for non-orphans, in which case the various algorithms presented earlier guarantee the transformed system is serially correct for all transaction names. The problem of ensuring that specific generic systems are serially correct for non-orphans is addressed elsewhere [4].

### 9.1.1. Generic Actions and Well-Formedness

For a generic system, we extend the object mapping as follows. Define  $\text{object}(\pi) = X$  if  $\pi$  is an  $\text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$  or  $\text{INFORM\_ABORT\_AT}(X)\text{OF}(T)$  action. We define the *generic actions* to be the serial actions, plus the  $\text{INFORM\_COMMIT}$  and  $\text{INFORM\_ABORT}$  actions.

Now we define "generic object well-formedness." Let  $X$  be any object name. A sequence  $\beta$  of generic actions  $\pi$  with  $\text{object}(\pi) = X$  is defined to be *generic object well-formed* for  $X$  provided that the following conditions hold.

1. There is at most one  $\text{CREATE}(T)$  event in  $\beta$  for any transaction  $T$ .
2. There is at most one  $\text{REQUEST\_COMMIT}$  event in  $\beta$  for any transaction  $T$ .
3. If there is a  $\text{REQUEST\_COMMIT}$  event for access transaction  $T$  in  $\beta$ , then there is a preceding  $\text{CREATE}(T)$  event in  $\beta$ .
4. There is no transaction  $T$  for which both an  $\text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$  event and an  $\text{INFORM\_ABORT\_AT}(X)\text{OF}(T)$  event occur.
5. If an  $\text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$  event occurs in  $\beta$  and  $T$  is an access to  $X$ , then there is a preceding  $\text{REQUEST\_COMMIT}$  event for  $T$ .

The following simple lemma shows a connection between transaction well-formedness for accesses and generic object well-formedness.

**Lemma 32:** Suppose  $\beta$  is a sequence of generic actions  $\pi$  with  $\text{object}(\pi) = X$ . If  $\beta$  is generic object well-formed for  $X$  and  $T$  is an access to  $X$ , then  $\beta|T$  is transaction well-formed for  $T$ .

### 9.1.2. Generic Object Automata

A *generic object automaton*  $G$  for an object name  $X$  is an I/O automaton with the following external action signature.

Input:

$\text{CREATE}(T)$ , for  $T$  an access to  $X$   
 $\text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$ , for  $T$  any transaction name  
 $\text{INFORM\_ABORT\_AT}(X)\text{OF}(T)$ , for  $T$  any transaction name

Output:

$\text{REQUEST\_COMMIT}(T,v)$ , for  $T$  an access to  $X$  and  $v$  a value

In addition,  $G$  may have an arbitrary set of internal actions.  $G$  is required to preserve generic object well-formedness. Except for this well-formedness requirement, generic object automata can be chosen arbitrarily.

Generic objects are similar to the abstract objects of Argus and other "object-oriented" systems. A generic object provides a set of accesses through which other transactions can observe and change the object's state. (These accesses can be thought of as instances of the "operations" usually assumed for objects in object-oriented systems.) Accesses can be invoked by concurrent transactions, and transactions can abort; thus, in generic transaction-processing systems that guarantee serial correctness, generic objects must provide synchronization and recovery. The objects studied in [11, 12], which use an exclusive locking variation of Moss's



algorithm [18] for synchronization combined with version stacks for recovery, are examples of generic objects that provide synchronization and recovery sufficient to ensure serial correctness for non-orphans. Similarly, the more general objects studied in [4], which use a commutativity-based locking algorithm that permits concurrent updates, are also generic objects that ensure serial correctness for non-orphans.

### 9.1.3. Generic Controller

The third kind of component in a generic system is the generic controller. The generic controller is also modeled as an I/O automaton. The transactions and generic objects have been specified to be any automata whose actions and behavior satisfy certain simple syntactic restrictions. A generic controller, however, is a fully specified automaton, particular to each system type. (Recall that we have assumed that the system type is fixed; we describe the generic controller for the fixed system type.)

The generic controller passes requests for the creation of subtransactions to the appropriate recipients, makes decisions about the commit or abort of transactions, passes reports about the completion of children back to their parents, and informs objects of the fate of transactions. It allows concurrency and aborts, and leaves the task of coping with them to the generic objects.

The generic controller is a very nondeterministic automaton. It may delay passing requests or reports or making decisions for arbitrary lengths of time, and may decide at any time to abort a transaction whose creation has been requested (but that has not yet completed). The generic controller can be implemented in many different ways by controllers that make specific choices from among the many nondeterministic possibilities. For instance, Moss [18] describes a distributed implementation of the generic controller that copes with node and communication failures yet still commits a subtransaction whenever possible. Our results apply *a fortiori* to all implementations of the generic controller obtained by restricting the nondeterminism.

The generic controller has the following action signature.

Input:

REQUEST\_CREATE( $T$ ),  $T \neq T_0$   
REQUEST\_COMMIT( $T, v$ )

Output:

CREATE( $T$ )  
COMMIT( $T$ ),  $T \neq T_0$   
ABORT( $T$ ),  $T \neq T_0$   
REPORT\_COMMIT( $T, v$ ),  $T \neq T_0$   
REPORT\_ABORT( $T$ ),  $T \neq T_0$   
INFORM\_COMMIT\_AT( $X$ )OF( $T$ ),  $T \neq T_0$   
INFORM\_ABORT\_AT( $X$ )OF( $T$ ),  $T \neq T_0$

The REQUEST\_CREATE and REQUEST\_COMMIT inputs are intended to be identified with the corresponding outputs of transaction and generic object automata, and correspondingly for the output actions.

Each state  $s$  of the generic controller consists of six sets:  $s.create\_requested$ ,  $s.created$ ,

$s.commit\_requested$ ,  $s.committed$ ,  $s.aborted$  and  $s.reported$ . The set  $s.commit\_requested$  is a set of (transaction,value) pairs, and the others are sets of transaction names. All are empty in the start state except for  $create\_requested$ , which is  $\{T_0\}$ . Define  $s.completed = s.committed \cup s.aborted$ .

The transition relation of the generic controller consists of exactly those triples  $(s',\pi,s)$  satisfying the preconditions and yielding the effects described below, where  $\pi$  is the indicated action. We include in the effects only those conditions on the state  $s$  that may change with the action. If a component of  $s$  is not mentioned in the effects, it is implicit that the set is the same in  $s'$  and  $s$ .

REQUEST\_CREATE( $T$ ),  $T \neq T_0$

Effect:

$s.create\_requested = s'.create\_requested \cup \{T\}$

REQUEST\_COMMIT( $T,v$ )

Effect:

$s.commit\_requested = s'.commit\_requested \cup \{(T,v)\}$

CREATE( $T$ )

Precondition:

$T \in s'.create\_requested - s'.created$

Effect:

$s.created = s'.created \cup \{T\}$

COMMIT( $T$ ),  $T \neq T_0$

Precondition:

$(T,v) \in s'.commit\_requested$  for some  $v$

$T \notin s'.completed$

Effect:

$s.committed = s'.committed \cup \{T\}$

ABORT( $T$ ),  $T \neq T_0$

Precondition:

$T \in s'.create\_requested - s'.completed$

Effect:

$s.aborted = s'.aborted \cup \{T\}$

REPORT\_COMMIT( $T,v$ ),  $T \neq T_0$

Precondition:

$T \in s'.committed$

$(T,v) \in s'.commit\_requested$

$T \notin s'.reported$

Effect:

$s.reported = s'.reported \cup \{T\}$

REPORT\_ABORT( $T$ ),  $T \neq T_0$

Precondition:

$T \in s'.aborted$

$T \notin s'.reported$

Effect:

$$s.\text{reported} = s'.\text{reported} \cup \{T\}$$

INFORM\_COMMIT\_AT(X)OF(T),  $T \neq T_0$

Precondition:

$$T \in s'.\text{committed}$$

INFORM\_ABORT\_AT(X)OF(T),  $T \neq T_0$

Precondition:

$$T \in s'.\text{aborted}$$

The generic controller assumes that its input actions, REQUEST\_CREATE and REQUEST\_COMMIT, can occur at any time, and simply records them in the appropriate components of the state. Once the creation of a transaction has been requested, the controller can create it by producing a CREATE action. The precondition of the CREATE action indicates that a given transaction will be created at most once; the effect of the CREATE is to record that the creation has occurred. Similarly, the effect of a COMMIT or ABORT action is to record that the action has occurred. REPORT\_COMMIT, REPORT\_ABORT, INFORM\_COMMIT and INFORM\_ABORT actions can be generated at any time after the corresponding COMMIT and ABORT actions have occurred. The precondition for a COMMIT action ensures that a transaction only commits if it has requested to do so, and it has not already completed (committed or aborted). The precondition for an ABORT action ensures that a transaction will be aborted only if a REQUEST\_CREATE has occurred for it and it has not already completed. There are no other constraints on when a transaction can be aborted, however. For example, a transaction can be aborted while some of its descendants are still running.

The following lemma follows easily by induction, using simple invariants maintained by the generic controller. (E.g., see [4].)

**Lemma 33:**

1. Let  $T$  be any transaction name. Then the generic controller preserves transaction well-formedness for  $T$ .
2. Let  $X$  be any object name. Then the generic controller preserves generic object well-formedness for  $X$ .
3. The generic controller preserves basic database well-formedness.

#### 9.1.4. Generic Database

A *generic database* is the composition of a strongly compatible set of automata indexed by the union of the set of object names and the singleton set  $\{GC\}$  (for "generic controller"). Associated with each object name  $X$  is a generic object automaton  $G_X$  for  $X$ , and associated with the name  $GC$  is the generic controller automaton for the system type.

**Lemma 34:** If  $\beta$  is a behavior of a generic database, then for every object name  $X$ ,  $\beta|X$  is generic object well-formed.

**Proof:** Let  $X$  be an object name. Of the components of the generic database, only  $G_X$  and the generic controller have external actions  $\pi$  for which  $\text{object}(\pi) = X$ . The other components thus trivially preserve generic object well-formedness for  $X$ . By explicit assumption,  $G_X$  preserves generic object well-formedness for  $X$ , and by Lemma 33, the generic controller preserves generic object well-formedness for  $X$ . The

result follows from Proposition 4. □

**Lemma 35:** Let  $B$  be a generic database. Then  $B$  preserves basic database well-formedness.

**Proof:** A simple case analysis. The only subtle case is when  $\beta\pi$  is a behavior of  $B$ , with  $\beta$  basic database well-formed and  $\pi$  a REQUEST\_COMMIT( $T, v$ ) event for an access  $T$  to  $X$ . The argument that  $\beta\pi|T$  is transaction well-formed for  $T$  depends upon the fact that  $\beta\pi|X$  is generic object well-formed for  $X$ , as shown in Lemma 34. □

It follows that the generic database is an example of a basic database.

### 9.1.5. Generic Systems

A *generic system* is the composition of a strongly compatible set of automata indexed by the union of the set of non-access transaction names, the set of object names and the singleton set  $\{GC\}$ . Associated with each non-access transaction name  $T$  is a transaction automaton  $A_T$  for  $T$ . Associated with each object name  $X$  is a generic object automaton  $G_X$  for  $X$ . Finally, associated with the name  $GC$  is the generic controller automaton for the system type.

When the particular generic system is clear from context, we call its executions, schedules and behaviors the *generic executions*, *generic schedules* and *generic behaviors*, respectively.

### 9.1.6. A Family of Affects Relations

Now we define a family  $R = \{R_\beta\}$  of affects relations for any particular generic system  $B$ . We do this by first defining a family  $R' = \{R'_\beta\}$  of directly-affects relations for  $B$ , and then taking transitive closures. For a sequence  $\beta$  of generic actions, define the relation  $R'_\beta$  to be the relation containing the pairs  $(\phi, \pi)$  of events such that  $\phi$  occurs before  $\pi$  in  $\beta$ , and at least one of the following holds:

- $\text{transaction}(\phi) = \text{transaction}(\pi)$  and  $\pi$  is an output event of the transaction,
- $\text{object}(\phi) = \text{object}(\pi)$  and  $\pi$  is a REQUEST\_COMMIT( $T, v$ ) event,
- $\phi$  is a REQUEST\_CREATE( $T$ ) and  $\pi$  a CREATE( $T$ ) event,
- $\phi$  is a REQUEST\_COMMIT( $T, v$ ) and  $\pi$  a COMMIT( $T$ ) event,
- $\phi$  is a REQUEST\_CREATE( $T$ ) and  $\pi$  an ABORT( $T$ ) event,
- $\phi$  is a COMMIT( $T$ ) and  $\pi$  a REPORT\_COMMIT( $T, v$ ) event,
- $\phi$  is an ABORT( $T$ ) and  $\pi$  a REPORT\_ABORT( $T$ ) event,
- $\phi$  is a COMMIT( $T$ ) and  $\pi$  an INFORM\_COMMIT\_AT( $X$ )OF( $T$ ) event, or
- $\phi$  is an ABORT( $T$ ) and  $\pi$  an INFORM\_ABORT\_AT( $X$ )OF( $T$ ) event.

Now define the relation  $R_\beta$  to be the transitive closure of  $R'_\beta$ . It is easy to see that  $R_\beta$  is an irreflexive partial order.

The idea is that  $\phi$  directly affects  $\pi$  if they both occur at the same transaction or object (and  $\pi$  is an output of the transaction, or a REQUEST\_COMMIT of the object), or if they involve different transactions or objects but the generic system will require  $\phi$  to occur before  $\pi$  can occur. This notion of one event affecting another is "safe," in the sense that  $\phi$  affects  $\pi$  if there is any

way that the precondition for  $\pi$  could require  $\phi$  to have occurred. If the events involve different transactions or objects, the preconditions for  $\pi$  in the generic controller require  $\phi$  to occur if  $\phi$  directly affects  $\pi$ . If the events occur at the same transaction or object, however, it might be that  $\phi$  happens to occur before  $\pi$ , yet that the particular transaction or object does not require  $\phi$  to occur before  $\pi$ . In the absence of more information about the particular transactions or objects used in a system, however, it is difficult to say more about the ways in which one event can affect another. Thus, we make the "safe" choice of assuming an effect whenever one could occur. Fortunately, the orphan management algorithms described earlier in this paper are essentially independent of the particular transactions and objects used in a system, and do not rely on more information about them.

The next few lemmas show that the family  $\{R_\beta\}$  defined above is a family of affects relations for B.

**Lemma 36:** If  $\beta$  is a finite behavior of generic system B and  $\gamma$  is an  $R_\beta$ -closed subsequence of  $\beta$ , then  $\gamma$  is a behavior of B.

**Proof:** For each non-access transaction name T, let  $A_T$  be the transaction automaton for T in B, and for each object name X, let  $G_X$  be the generic object automaton for X in B. By Proposition 2, it suffices to show that  $\gamma|T$  is a behavior of  $A_T$  for all non-access transaction names T, that  $\gamma|X$  is a behavior of  $G_X$  for all object names X, and that  $\gamma$  is a behavior of the generic controller. We show these in turn.

First, suppose that T is a non-access transaction name. If  $\gamma|T$  contains no output events of  $A_T$ , then the input-enabling property implies that  $\gamma|T$  is a behavior of  $A_T$ . So assume that there is at least one output event of  $A_T$  in  $\gamma|T$ , and let  $\pi$  be the last such event. Let  $\gamma'$  be the prefix of  $\gamma$  ending with  $\pi$ . Since  $\gamma$  is closed in  $\beta$ , it follows from the definition of  $R_\beta$  that  $\gamma$  contains all events of T that precede  $\pi$  in  $\beta$ . Thus,  $\gamma'|T$  is a prefix of  $\beta|T$  and so is a behavior of  $A_T$ . Since  $\gamma|T$  differs from  $\gamma'|T$  only by the possible inclusion of some final input events of  $A_T$ , the input-enabling property implies that  $\gamma|T$  is a behavior of  $A_T$ .

A similar argument shows that if X is an object name, then  $\gamma|X$  is a behavior of  $G_X$ .

Now we show that  $\gamma$  is a behavior of the generic controller. Note that the generic controller is deterministic in the sense that for a given state  $s'$  and action  $\pi$ , there is at most one state  $s$  such that  $(s', \pi, s)$  is a step of the generic controller. We proceed by induction on the lengths of prefixes  $\delta$  of  $\gamma$ . The basis, where the length of  $\delta$  is 0, is obvious. So suppose that  $\delta = \delta'\pi$ , where  $\pi$  is a single event. Let  $\beta'\pi$  be the prefix of  $\beta$  ending with  $\pi$ . We consider cases, showing in each case that  $\pi$  is enabled after  $\delta'$ .

1.  $\pi = \text{CREATE}(T)$

Then Lemma 35 and Proposition 5 imply that  $\text{REQUEST\_CREATE}(T)$  occurs in  $\beta'$ , and no  $\text{CREATE}(T)$  occurs in  $\beta'$ . Since  $\gamma$  is  $R_\beta$ -closed in  $\beta$ ,  $\text{REQUEST\_CREATE}(T)$  also occurs in  $\delta'$ . Similarly, no  $\text{CREATE}(T)$  occurs in  $\delta'$ . It follows that  $\pi$  is enabled in the (unique) state in which  $\delta'$  can leave the generic controller. Thus,  $\pi$  is enabled after  $\delta'$ .

2.  $\pi = \text{COMMIT}(T)$ 

Then Lemma 35 and Proposition 5 imply that  $\beta'$  contains  $\text{REQUEST\_COMMIT}(T, v)$  and contains no completion events for  $T$ . Since  $\gamma$  is  $R_\beta$ -closed in  $\beta$ ,  $\delta'$  contains  $\text{REQUEST\_COMMIT}(T, v)$ , and does not contain a completion event for  $T$ . It follows that  $\pi$  is enabled after  $\delta'$ .

3.  $\pi = \text{ABORT}(T)$ 

Then  $\beta'$  contains  $\text{REQUEST\_CREATE}(T)$  and contains no completion events for  $T$ , so  $\delta'$  contains  $\text{REQUEST\_CREATE}(T)$  and no completion events for  $T$ . Thus,  $\pi$  is enabled after  $\delta'$ .

4.  $\pi = \text{REPORT\_COMMIT}(T, v)$ 

Then  $\beta'$  contains  $\text{COMMIT}(T)$  and  $\text{REQUEST\_COMMIT}(T, v)$ , and contains no report events for  $T$ . Therefore, the same is true of  $\delta'$ , so  $\pi$  is enabled after  $\delta'$ .

5.  $\pi = \text{REPORT\_ABORT}(T)$ 

Similar to the preceding arguments.

6.  $\pi = \text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$ 

Similar to the preceding arguments.

7.  $\pi = \text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$ 

Similar to the preceding arguments.

□

**Lemma 37:** The family  $\{R_\beta\}$  is a family of affects relations for  $B$ .

**Proof:** The first two properties of families of affects relations are immediate from the definition of  $\{R_\beta\}$ . The third property is proven in the lemma above.

The fourth property requires that whenever an  $\text{ABORT}(T)$  event  $\phi$  affects an event  $\pi$  of certain types, there is a specific type of intervening event  $\psi$  that is also affected by  $\phi$ . To see that this property is satisfied, note that  $R_\beta$  is defined as the transitive closure of the directly-affects relation  $R'_\beta$ . Each type of event  $\pi$  of interest is only directly affected by earlier events that satisfy the restrictions on  $\psi$ . Thus, if  $\pi$  is affected by  $\phi$ , it must be so affected by the transitive closure over a chain of directly-affects relations in which an event  $\psi$  of the appropriate type occurs. □

### 9.1.7. Applying Orphan Management Algorithms to Generic Systems

Since we have shown that generic systems are instances of basic systems, and that the directly-affects relation we defined above generates a family of affects relations, we may apply our general results for orphan management algorithms to generic systems. Before we state these results, we illustrate one orphan management transformation by describing explicitly the steps of the system obtained by applying the Argus algorithm to a generic system.

If  $B$  is a generic system and  $R'_\beta$  is the family of directly-affects relations given above, then the steps of  $\text{Argus}(B, R'_\beta)$  are defined as follows.

- $(s'.\text{basic\_state}, \pi, s.\text{basic\_state})$  is a step of the generic database.
- $s.\text{history} = s'.\text{history}\pi$  if  $\pi$  is a basic action.
- $s.\text{history} = s'.\text{history}$  if  $\pi$  is not a basic action.

- If  $\phi \neq \pi$ , then  $s.\text{known\_aborts}(\phi) = s'.\text{known\_aborts}(\phi)$ .
- If  $\pi$  is a  $\text{REQUEST\_CREATE}(T)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\phi$  in  $s'.\text{history}$  such that  $\text{transaction}(\phi) = \text{parent}(T)$ .
- If  $\pi$  is a  $\text{REQUEST\_COMMIT}(T,v)$  action, where  $T$  is a non-access transaction name, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\phi$  in  $s'.\text{history}$  such that  $\text{transaction}(\phi) = T$ .
- If  $\pi$  is a  $\text{REQUEST\_COMMIT}(T,v)$  action, where  $T$  is an access transaction name, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\phi$  in  $s'.\text{history}$  such that  $\text{object}(\phi) = \text{object}(T)$ .
- If  $\pi$  is a  $\text{CREATE}(T)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\text{REQUEST\_CREATE}(T)$  events  $\phi$  in  $s'.\text{history}$ .
- If  $\pi$  is a  $\text{COMMIT}(T)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\text{REQUEST\_COMMIT}(T,v)$  events  $\phi$  in  $s'.\text{history}$ .
- If  $\pi$  is an  $\text{ABORT}(T)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\text{REQUEST\_CREATE}(T)$  events  $\phi$  in  $s'.\text{history}$ .
- If  $\pi$  is a  $\text{REPORT\_COMMIT}(T,v)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\text{COMMIT}(T)$  events  $\phi$  in  $s'.\text{history}$ .
- If  $\pi$  is a  $\text{REPORT\_ABORT}(T)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\text{ABORT}(T)$  events  $\phi$  in  $s'.\text{history}$ , and  $T \in s.\text{known\_aborts}(\pi)$ .
- If  $\pi$  is an  $\text{INFORM\_COMMIT\_AT}(X)\text{OF}(T)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\text{COMMIT}(T)$  events  $\phi$  in  $s'.\text{history}$ .
- If  $\pi$  is an  $\text{INFORM\_ABORT\_AT}(X)\text{OF}(T)$  action, then  $s'.\text{known\_aborts}(\phi) \subseteq s.\text{known\_aborts}(\pi)$  for all  $\text{ABORT}(T)$  events  $\phi$  in  $s'.\text{history}$ , and  $T \in s.\text{known\_aborts}(\pi)$ .
- If  $\pi$  is a  $\text{REQUEST\_COMMIT}(T,v)$  action for an access  $T$  to object  $X$ , then there is no ancestor of  $T$  in  $s'.\text{known\_aborts}(\phi)$ , for any event  $\phi$  of object  $X$  in  $s'.\text{history}$ .

The filtered database uses global information about the history to prevent the  $\text{REQUEST\_COMMIT}$  of an access from occurring if it would be affected by an  $\text{ABORT}$  of an ancestor. The Argus database uses more local information, which is obtained by propagating the  $\text{known\_aborts}$  sets from each event to any later events that it directly affects. For example, consider the application above of the Argus algorithm to a generic database. The  $\text{known\_aborts}$  set for a  $\text{REQUEST\_CREATE}(T)$  action is obtained from the  $\text{known\_aborts}$  sets for all preceding events at  $\text{parent}(T)$ . Since  $\text{REQUEST\_CREATE}(T)$  is generated by  $\text{parent}(T)$ , the  $\text{known\_aborts}$  set for it can easily be computed with information available locally at  $\text{parent}(T)$  when the  $\text{REQUEST\_CREATE}(T)$  action occurs. A similar situation arises with  $\text{REQUEST\_COMMIT}$  actions, which are outputs of transactions and objects. The other actions, which are outputs of the generic controller, are directly affected by exactly one preceding event. Thus, the  $\text{known\_aborts}$  set for one of these actions can easily be computed from the  $\text{known\_aborts}$  set for the single event that precedes it. For instance, the  $\text{known\_aborts}$  set for a  $\text{CREATE}(T)$  event can be obtained directly from the  $\text{known\_aborts}$  set for the preceding  $\text{REQUEST\_CREATE}(T)$  event. If the two events occur at the same site in a network, this

information would be available locally; if they occur at different sites, it could be sent in the message used to transmit the REQUEST\_CREATE event to the site that performs the CREATE event.

The following corollary shows that the Argus algorithm and the clocked algorithm from [16] can both be used for a generic system:

**Corollary 38:** Let  $B$  be a generic system, and  $R$  and  $R'$  the family of affects and directly-affects relations defined above. If  $B$  is serially correct for non-orphans, then the following are true:

- Filtered( $B, R$ ) is serially correct.
- Argus( $B, R'$ ) is serially correct.
- Strictly-Filtered( $B$ ) is serially correct.
- Clocked( $B$ ) is serially correct.

## 9.2. Pseudotime Systems

The essential feature of systems using timestamps is the explicit construction of a sibling order representing the intended serialization of an execution. This order is represented in terms of intervals of *pseudotime*, an arbitrarily chosen totally ordered set. Formally, we let  $\mathcal{P}$  be the set of pseudotimes, ordered by  $<$ . We represent pseudotime intervals as half-open intervals  $[p, q)$  in  $\mathcal{P}$ , and refer to them using capital letters. If  $P = [p, q)$ , then we write  $P_{\min}$  for  $p$  and  $P_{\max}$  for  $q$ . If  $P$  and  $Q$  are intervals of pseudotime, we write  $P < Q$  if  $P_{\max} \leq Q_{\min}$ . Clearly, if  $P < Q$ , then  $P$  and  $Q$  are disjoint.

A pseudotime system consists of a collection of transaction automata, one for each non-access transaction name, a collection of "pseudotime object automata," one for each object name, and a single "pseudotime controller automaton." The interactions between the components are as follows. The transaction interface is exactly as before. A pseudotime object automaton for  $X$  has the same actions as a generic object automaton, with the addition of INFORM\_TIME\_AT( $X$ )OF( $T, p$ ) input actions to inform the object that pseudotime  $p$  has been assigned to an access transaction  $T$ . The pseudotime controller has the same actions as the generic controller, with the addition of INFORM\_TIME\_AT( $X$ )OF( $T, p$ ) output actions (for access transactions  $T$ ) and ASSIGN\_PSEUDOTIME( $T, P$ ) output actions (for all transactions  $T$ ) by which the controller assigns the pseudotime range  $P$  to transaction  $T$ .

### 9.2.1. Pseudotime Actions and Well-formedness

The object mapping for a pseudotime system is the same as that for a generic system, with the addition that we define  $\text{object}(\pi) = X$  if  $\pi$  is an INFORM\_TIME\_AT( $X$ )OF( $T, p$ ) action. We define the *pseudotime actions* to be the generic actions, plus the INFORM\_TIME and ASSIGN\_PSEUDOTIME actions.

Now we define "pseudotime object well-formedness." Let  $X$  be any object name. A sequence of pseudotime actions  $\pi$  with  $\text{object}(\pi) = X$  is defined to be *pseudotime object well-formed* for  $X$  provided that the following conditions hold.

1. There is at most one CREATE( $T$ ) event in  $\beta$  for any transaction  $T$ .



2. There is at most one REQUEST\_COMMIT event in  $\beta$  for any transaction  $T$ .
3. If there is a REQUEST\_COMMIT event for  $T$  in  $\beta$ , then there is a preceding CREATE( $T$ ) event and also a preceding INFORM\_TIME\_AT( $X$ )OF( $T, p$ ) in  $\beta$ .
4. There is no transaction  $T$  for which there are two different pseudotimes,  $p$  and  $p'$ , such that INFORM\_TIME\_AT( $X$ )OF( $T, p$ ) and INFORM\_TIME\_AT( $X$ )OF( $T, p'$ ) both occur in  $\beta$ .
5. There is no pseudotime  $p$  for which there are two different transactions,  $T$  and  $T'$ , such that INFORM\_TIME\_AT( $X$ )OF( $T, p$ ) and INFORM\_TIME\_AT( $X$ )OF( $T', p$ ) both occur in  $\beta$ .
6. There is no transaction  $T$  for which both an INFORM\_COMMIT\_AT( $X$ )OF( $T$ ) event and an INFORM\_ABORT\_AT( $X$ )OF( $T$ ) event occur.
7. If an INFORM\_COMMIT\_AT( $X$ )OF( $T$ ) event occurs in  $\beta$  and  $T$  is an access to  $X$ , then there is a preceding REQUEST\_COMMIT event for  $T$ .

### 9.2.2. Pseudotime Object Automata

A *pseudotime object automaton*  $P$  for an object name  $X$  is an I/O automaton with the following external action signature.

Input:

CREATE( $T$ ),  $T$  an access to  $X$   
 INFORM\_COMMIT\_AT( $X$ )OF( $T$ )  
 INFORM\_ABORT\_AT( $X$ )OF( $T$ )  
 INFORM\_TIME\_AT( $X$ )OF( $T, p$ ),  $T$  an access to  $X$ ,  $p \in P$

Output:

REQUEST\_COMMIT( $T, v$ ),  $T$  an access to  $X$

In addition,  $P$  may have an arbitrary set of internal actions.  $P$  is required to preserve pseudotime object well-formedness.

### 9.2.3. Pseudotime Controller

The *pseudotime controller* guarantees that siblings are assigned disjoint intervals of pseudotime, and that each transaction's interval is a subset of that of its parent. The pseudotime controller has the actions of the generic controller together with an extra class of output actions ASSIGN\_PSEUDOTIME( $T, P$ ) for  $T \neq T_0$  and  $P$  a pseudotime interval. The purpose of the ASSIGN\_PSEUDOTIME actions is to construct, at run-time, a sibling order that specifies the apparent serial ordering of transactions. Also, there is an extra class of actions INFORM\_TIME\_AT( $X$ )OF( $T, p$ ) for access transactions  $T$ . A state  $s$  of the pseudotime controller has the same components as a state of the generic controller together with an additional component  $s.\text{interval}$ , which is a partial function from  $\mathcal{T}$  to the set of pseudotime intervals. In the initial state  $s_0$  of the pseudotime controller  $s_0.\text{interval} = \{(T_0, P_0)\}$  for some pseudotime interval  $P_0$ , and all other components are as in the initial state of the generic controller.

The transition relation for generic actions is the same as that for the generic controller, except that the actions CREATE( $T$ ) and ABORT( $T$ ) have an additional precondition:  $T \in$

$\text{domain}(s'.\text{interval})$ . The additional actions are determined as follows.

**ASSIGN\_PSEUDOTIME(T,P)**

Precondition:

$T \in s'.\text{create-requested}$

$T \notin \text{domain}(s'.\text{interval})$

$P \subseteq s'.\text{interval}(\text{parent}(T))$

$P > s'.\text{interval}(T')$  for every  $T'$  in  $\text{siblings}(T) \cap \text{domain}(s'.\text{interval})$

Effect:

$s'.\text{interval} = s'.\text{interval} \cup \{(T,P)\}$

**INFORM\_TIME\_AT(X)OF(T,p)**,  $T$  an access to  $X$

Precondition:

$(T,P) \in s'.\text{interval}$

$p = P_{\min}$

The following lemma is straightforward.

**Lemma 39:**

1. Let  $T$  be any transaction name. Then the pseudotime controller preserves transaction well-formedness for  $T$ .
2. Let  $X$  be any object name. Then the pseudotime controller preserves pseudotime object well-formedness for  $X$ .
3. The pseudotime controller preserves basic database well-formedness.

#### 9.2.4. Pseudotime Database

A *pseudotime database* is the composition of a strongly compatible set of automata indexed by the union of the set of object names and the singleton set  $\{PC\}$  (for "pseudotime controller"). Associated with each object name  $X$  is a pseudotime object automaton  $P_X$  for  $X$ . Finally, associated with the name  $PC$  is the pseudotime controller automaton for the system type.

**Lemma 40:** If  $\beta$  is a behavior of a pseudotime database, then for every object name  $X$ ,  $\beta|X$  is pseudotime object well-formed.

**Lemma 41:** Let  $B$  be a pseudotime database. Then  $B$  preserves basic database well-formedness.

It follows that the pseudotime database is an example of a basic database.

#### 9.2.5. Pseudotime Systems

A *pseudotime system* is the composition of a strongly compatible set of automata indexed by the union of the set of non-access transaction names, the set of object names and the singleton set  $\{PC\}$ . Associated with each non-access transaction name  $T$  is a transaction automaton  $A_T$  for  $T$ . Associated with each object name  $X$  is a pseudotime object automaton  $P_X$  for  $X$ . Finally, associated with the name  $PC$  is the pseudotime controller automaton for the system type.

When the particular pseudotime system is clear from context, we call its executions, schedules and behaviors the *pseudotime executions*, *pseudotime schedules* and *pseudotime behaviors*, respectively.

### 9.2.6. A Family of Affects Relations

Now we define a family  $R = \{R_\beta\}$  of affects relations for any particular pseudotime system  $B$ . We do this by first defining a family  $R' = \{R'_\beta\}$  of directly-affects relations for  $B$ , and then taking transitive closures. For a sequence  $\beta$  of pseudotime actions, define the relation  $R'_\beta$  to be the relation containing the pairs  $(\phi, \pi)$  of events such that  $\phi$  occurs before  $\pi$  in  $\beta$ , and at least one of the following holds:

- $\text{transaction}(\phi) = \text{transaction}(\pi)$  and  $\pi$  is an output event of the transaction,
- $\text{object}(\phi) = \text{object}(\pi)$  and  $\pi$  is a `REQUEST_COMMIT(T,v)` event,
- $\phi$  is a `REQUEST_CREATE(T)` and  $\pi$  an `ASSIGN_PSEUDOTIME(T,P)` event,
- $\phi$  is an `ASSIGN_PSEUDOTIME(T,P)` and  $\pi$  a `CREATE(T)` event,
- $\phi$  is a `REQUEST_COMMIT(T,v)` and  $\pi$  a `COMMIT(T)` event,
- $\phi$  is a `REQUEST_CREATE(T)` and  $\pi$  an `ABORT(T)` event,
- $\phi$  is an `ASSIGN_PSEUDOTIME(T,P)` and  $\pi$  an `ABORT(T)` event,
- $\phi$  is a `COMMIT(T)` and  $\pi$  a `REPORT_COMMIT(T,v)` event,
- $\phi$  is an `ABORT(T)` and  $\pi$  a `REPORT_ABORT(T)` event,
- $\phi$  is a `COMMIT(T)` and  $\pi$  an `INFORM_COMMIT_AT(X)OF(T)` event,
- $\phi$  is an `ABORT(T)` and  $\pi$  an `INFORM_ABORT_AT(X)OF(T)` event, or
- $\phi$  is an `ASSIGN_PSEUDOTIME(T,P)` event and  $\pi$  an `INFORM_TIME_AT(X)OF(T,p)` event.

Once again, define the relation  $R_\beta$  to be the transitive closure of  $R'_\beta$ . It is easy to see that  $R_\beta$  is an irreflexive partial order. We claim that the family  $\{R_\beta\}$  defined above is a family of affects relations for  $G$ .

**Lemma 42:** If  $\beta$  is a behavior of pseudotime system  $B$  and  $\gamma$  is an  $R_\beta$ -closed subsequence of  $\beta$ , then  $\gamma$  is a behavior of  $B$ .

**Proof:** Analogous to the proof of Lemma 36. □

**Lemma 43:** The family  $\{R_\beta\}$  is a family of affects relations for  $B$ .

### 9.2.7. Applying Orphan Management Algorithms to Pseudotime Systems

The following easy corollary shows that the Argus algorithm and the clocked algorithm from [16] can both be used for a pseudotime system:

**Corollary 44:** Let  $B$  be a pseudotime system, and  $R$  and  $R'$  the family of affects and directly-affects relations defined above. If  $B$  is serially correct for non-orphans, then the following are true:

- $\text{Filtered}(B,R)$  is serially correct.
- $\text{Argus}(B,R')$  is serially correct.
- $\text{Strictly-Filtered}(B)$  is serially correct.
- $\text{Clocked}(B)$  is serially correct.

## 10. Conclusions

We have defined correctness properties for orphan management algorithms, and have presented precise descriptions and proofs for two algorithms from [9] and [16]. Our proofs are quite simple, and show that the systems exhibit a substantial degree of modularity: the orphan management algorithms can be used in combination with any concurrency control protocol (in basic system form) that is serially correct for non-orphans. The simplicity of our proofs is a direct result of this modularity, and is in sharp contrast to earlier work [6], in which the orphan management algorithm and the concurrency control protocol were not cleanly separated.

Our proofs have an interesting structure. We first define a simple abstract algorithm that uses global information about the history of the system, and show that it ensures that orphans see consistent views. We then formalize the Argus algorithm and the clocked algorithm in a way that only requires the use of local information, and show that each simulates the more abstract algorithm. The simulation proofs are quite simple, and do not require re-proving the properties already proved for the abstract algorithm. The correctness of the Argus and clocked algorithms then follows directly from the correctness of the abstract algorithm.

In this paper we have analyzed only orphans that result from aborts of transactions. Interesting algorithms have also been developed for detecting and eliminating orphans arising from crashes [9, 16]. These algorithms seem more complicated than the algorithms for handling aborts. An open question is whether the known algorithms for handling crash orphans can be analyzed using techniques similar to those in this paper. In particular, it would be nice to find a similar separation of concerns for those algorithms, so that the crash-orphan algorithms can be understood independently of concurrency control protocols and abort-orphan algorithms. Whether this will be possible is still unknown.

## 11. Acknowledgements

We thank Alan Fekete, Ken Goldman and Sharon Perl for their comments on earlier versions of this work.

## References

- [1] Allchin, J. E.  
*An architecture for reliable decentralized systems.*  
PhD thesis, Georgia Institute of Technology, September, 1983.  
Available as Technical Report GIT-ICS-83/23.
- [2] Aspnes, J., Fekete, A., Lynch, N., Merritt, M. and Weihl, W.  
A Theory of Timestamp-Based Concurrency Control for Nested Transactions.  
In *Proceedings of 14th International Conference on Very Large Data Bases*, pages 431-444. August, 1988.
- [3] Bernstein, P., Hadzilacos, V. and Goodman, N.  
*Concurrency Control and Recovery in Database Systems.*  
Addison-Wesley, 1987.
- [4] Fekete, A., Lynch, N., Merritt, M. and Weihl, W.  
*Commutativity-Based Locking for Nested Transactions.*  
Technical Memo MIT/LCS/TM-370, Massachusetts Institute Technology, Laboratory for Computer Science, August, 1988.  
Revised version to appear in JCSS.
- [5] Goldman, K. and Lynch, N.  
Nested Transactions and Quorum Consensus.  
In *Proceedings of 6th ACM Symposium on Principles of Distributed Computation*, pages 27-41. August, 1987.  
Expanded version available as Technical Report MIT/LCS/TM-390, Laboratory for Computer Science, Massachusetts Institute Technology, Cambridge, MA, May 1987.
- [6] Goree, J. A.  
Internal consistency of a distributed transaction system with orphan detection.  
Master's thesis, Massachusetts Institute of Technology, January, 1983.  
Available as MIT/LCS/TR-286.
- [7] D.L. Detlefs, M.P. Herlihy, and J. M. Wing.  
Inheritance of Synchronization and Recovery Properties in Avalon/C++.  
*IEEE Computer*, December, 1988.
- [8] Liskov, B., and Scheifler, R.  
Guardians and actions: linguistic support for robust, distributed programs.  
*ACM Transactions on Programming Languages and Systems* 5(3):381-404, July, 1983.
- [9] Liskov, B., Scheifler, R., Walker, E. F., and Weihl, W.  
Orphan Detection.  
In *Proceedings of the Seventeenth International Symposium on Fault-Tolerant Computing*. IEEE, July, 1987.
- [10] Lynch, N. A.  
Concurrency control for resilient nested transactions.  
*Advances in Computing Research* 3:335-373, 1986.

- [11] Lynch, N. A., and Merritt, M.  
*Introduction to the theory of nested transactions.*  
Technical Report MIT-LCS-TR-367, Massachusetts Institute of Technology, 1986.  
Also in *Theoretical Computer Science* 62 (1988), pages 123-185.
- [12] Lynch, N. and Merritt, M.  
Introduction to the Theory of Nested Transactions.  
In *International Conference on Database Theory*, pages 278-305. Rome, Italy, September, 1986.
- [13] Lynch, N. and Merritt, M. and Weihl, W. and Fekete, A.  
*Atomic Transactions.*  
In preparation.
- [14] Lynch, N. and Tuttle, M.  
Hierarchical Correctness Proofs for Distributed Algorithms.  
In *Proceedings of 6th ACM Symposium on Principles of Distributed Computing*, pages 137-151. August, 1987.  
Expanded version available as Technical Report MIT/LCS/TR-387, Laboratory for Computer Science, Massachusetts Institute Technology, Cambridge, MA., April 1987.
- [15] Lynch, N. and Tuttle, M.  
An Introduction to Input/Output Automata.  
To be published in *Centrum voor Wiskunde en Informatica Quarterly*. Also in Technical Memo, MIT/LCS/TM-373, Lab for Computer Science Massachusetts Institute of Technology, November 1988.
- [16] McKendry, M., and Herlihy, M.  
Time-driven orphan elimination.  
In *Proceedings of the 5th Symposium on Reliability in Distributed Software and Database Systems*, pages 42-48. IEEE, January, 1986.
- [17] McKendry, M., and Herlihy, M.  
*Timestamp-based orphan elimination.*  
Technical Report CMU-CS-87-108, Carnegie-Mellon University, 1987.
- [18] Moss, J. E. B.  
*Nested transactions: an approach to reliable distributed computing.*  
PhD thesis, Massachusetts Institute of Technology, 1981.  
Available as Technical Report MIT/LCS/TR-260.
- [19] Nelson, B. J.  
*Remote procedure call.*  
PhD thesis, Carnegie-Mellon University Department of Computer Science, May, 1981.  
Available as CMU-CS-81-119.
- [20] Perl, S.  
Distributed Commit Protocols for Nested Atomic Actions.  
Master's thesis, Massachusetts Institute Technology, September, 1987.  
Available as MIT/LCS/TR-431.

- [21] Pu, C., and Noe, J. D.  
*Nested transactions for general objects: the Eden implementation.*  
Technical Report TR-85-12-03, University of Washington Department of Computer Science, December, 1985.
- [22] Rosenkrantz, D. J., Lewis, P. M., and Stearns, R. E.  
System Level Concurrency Control for Distributed Database Systems.  
*ACM Transactions on Database Systems* 3(2):178-198, June, 1978.
- [23] Spector, A. and Swedlow, K.  
Guide to the Camelot Distributed Transaction Facility: Release 1.  
October, 1987  
Available from Carnegie Mellon University, Pittsburgh, PA.
- [24] Walker, E. F.  
Orphan Detection in the Argus System.  
Master's thesis, Massachusetts Institute of Technology, May, 1984.  
Available as MIT/LCS/TR-326.
- [25] Weihl, W. E.  
*Specification and implementation of atomic data types.*  
PhD thesis, Massachusetts Institute of Technology, 1984.  
Available as Technical Report MIT/LCS/TR-314.
- [26] Weihl, W. E.  
Local atomicity properties: modular concurrency control for abstract data types.  
*ACM Transactions on Programming Languages and Systems* , April, 1989.

## Table of Contents

<b>1. Introduction</b>	<b>1</b>
1.1. New Results	1
1.2. Related Work	2
1.3. Organization of this Paper	2
<b>2. Formal Preliminaries</b>	<b>4</b>
2.1. The Input/Output Automaton Model	4
2.1.1. Action Signatures	5
2.1.2. Input/Output Automata	5
2.1.3. Executions, Schedules and Behaviors	6
2.2. Composition	7
2.2.1. Composition of Action Signatures	7
2.2.2. Composition of Automata	8
2.2.3. Properties of Systems of Automata	8
2.3. Implementation	9
2.4. Preserving Properties	10
<b>3. Basic Systems</b>	<b>10</b>
3.1. Overview	11
3.2. System Types	11
3.3. General Structure of Basic Systems	12
3.4. Serial Actions	15
3.4.1. Terminology	15
3.4.2. Well-Formedness	16
3.5. Basic Systems	17
3.5.1. Transaction Automata	17
3.5.2. Basic Database Automata	18
3.5.3. Basic Systems	19
3.6. Serial Correctness	20
<b>4. Information Flow</b>	<b>21</b>
4.1. Families of Affects Relations	21
4.2. Families of Directly-Affects Relations	22
4.3. Using Families of Affects Relations to Describe Orphan Management Algorithms	23
<b>5. Filtered Systems</b>	<b>23</b>
5.1. The Filtered Database	24
5.2. The Filtered System	25
5.3. Simulation of the Basic System by the Filtered System	26
<b>6. Argus Systems</b>	<b>27</b>
6.1. The Argus Database	27
6.2. The Argus System	29
6.3. Simulation of the Basic System by the Argus System	30
<b>7. Strictly Filtered Systems</b>	<b>31</b>
7.1. The Strictly Filtered Database	32
7.2. The Strictly Filtered System	32
7.3. Simulation of the Basic System by the Strictly Filtered System	32
<b>8. Clocked Systems</b>	<b>33</b>
8.1. The Clocked Database	33
8.2. The Clocked System	35
8.3. Simulation of the Basic System by the Clocked System	35
<b>9. Examples</b>	<b>36</b>
9.1. Generic Systems	36



9.1.1. Generic Actions and Well-Formedness	37
9.1.2. Generic Object Automata	37
9.1.3. Generic Controller	38
9.1.4. Generic Database	40
9.1.5. Generic Systems	41
9.1.6. A Family of Affects Relations	41
9.1.7. Applying Orphan Management Algorithms to Generic Systems	43
9.2. Pseudotime Systems	45
9.2.1. Pseudotime Actions and Well-formedness	45
9.2.2. Pseudotime Object Automata	46
9.2.3. Pseudotime Controller	46
9.2.4. Pseudotime Database	47
9.2.5. Pseudotime Systems	47
9.2.6. A Family of Affects Relations	48
9.2.7. Applying Orphan Management Algorithms to Pseudotime Systems	48
10. Conclusions	49
11. Acknowledgements	49

OFFICIAL DISTRIBUTION LIST

Director 2 copies  
Information Processing Techniques Office  
Defense Advanced Research Projects Agency  
1400 Wilson Boulevard  
Arlington, VA 22209

Office of Naval Research 2 copies  
800 North Quincy Street  
Arlington, VA 22217  
Attn: Dr. Gary Koop, Code 433

Director, Code 2627 6 copies  
Naval Research Laboratory  
Washington, DC 20375

Defense Technical Information Center 12 copies  
Cameron Station  
Alexandria, VA 22314

National Science Foundation 2 copies  
Office of Computing Activities  
1800 G. Street, N.W.  
Washington, DC 20550  
Attn: Program Director

Dr. E.B. Royce, Code 38 1 copy  
Head, Research Department  
Naval Weapons Center  
China Lake, CA 93555